# <u>תרגיל בית מספר 6 - להגשה עד 10/01/2022 בש</u>עה 23:55

קראו בעיון את הנחיות העבודה <u>וההגשה</u> המופיעות באתר הקורס, תחת התיקייה assignments. חריגה מההנחיות תגרור ירידת ציון / פסילת התרגיל.

#### : הגשה

- תשובותיכם יוגשו בקובץ pdf ובקובץ py בהתאם להנחיות בכל שאלה.
- השתמשו בקובץ השלד skeleton6.py כבסיס לקובץ ה py אותו אתם מגישים. לא לשכוח לשנות את שם הקובץ למספר ת"ז שלכם לפני ההגשה, עם סיומת py.
- בסהייכ מגישים שני קבצים בלבד. עבור סטודנטית שמספר תייז שלה הוא 012345678 הקבצים שיש להגיש הם  $hw6_012345678.pyf$ .
  - הקפידו לענות על כל מה שנשאלתם.
  - תשובות מילוליות והסברים צריכים להיות תמציתיים, קולעים וברורים.
     להנחיה זו מטרה כפולה:
    - 1. על מנת שנוכל לבדוק את התרגילים שלכם בזמן סביר.
- 2. כדי להרגיל אתכם להבעת טיעונים באופן מתומצת ויעיל, ללא פרטים חסרים מצד אחד אך ללא עודף בלתי הכרחי מצד שני. זוהי פרקטיקה חשובה במדעי המחשב.

## שאלה 1 (30 נקי)

שאלה זו עוסקת בדקדוקים חסרי הקשר.

תזכרות: בהינתן אלפבית  $\Sigma$ , שפה היא קבוצה (אולי אינסופית) של מחרוזות מעל האלפבית  $\Sigma$ , כלומר קבוצה של מחרוזות כך שכל מחרוזת מכילה תווים רק מתוך הקבוצה  $\Sigma$ .

לדוגמה אם  $\Sigma=\{0,1,2\}$  אז  $\Sigma=\{0,1,2,22\}$  היא שפה (סופית) מעל  $\Sigma$ . בהינתן שפה  $L=\{\varepsilon,0,01,02,22\}$  אז בהיא שפה חסרת  $C=\{0,1,2\}$  בהיעתן לגזור באמצעות הדקדוק היא הקשר אם קיים דקדוק חסר הקשר  $C=\{0,1,2\}$  הדקדוק היא העפה של הדקדוק  $C=\{0,1,2\}$ , ונסמן ונסמן  $C=\{0,1,2\}$  במקרה זה נאמר ש- $C=\{0,1,2\}$  היא השפה של הדקדוק  $C=\{0,1,2\}$ , ונסמן  $C=\{0,1,2\}$ 

א. עבור כל אחת מהשפות הבאות, הראו שהשפה היא שפה חסרת הקשר על ידי כתיבת דקדוק מתאים. הגדירו באופן א. עבור כל אחת מהשפות הבאות, הראו שהשפה היא שפה חסרת הקשר אור צורך להוכיח את פורמלי את כל אחד מן הפרמטרים  $V, \Sigma, R$  (השתמשו באות S בתור משתנה הפתרון שלכם. הדקדוקים אינם צריכים להיות בצורת CNF.

**.0 היא** היא קבוצת המספרים הטבעיים כולל  $\mathbb{N}$  היא קבוצת המספרים הטבעיים

- $L_1 = \{0^n 12^{2n} : n \in \mathbb{N}\}$  .i
- ופעולות איים +,\* וכן עם סוגריים חוקיים .ii השפה ביטויים אל כל הביטויים המתמטיים שניתן לייצר עם x,y,z ופעולות הביטויים המתמטיים חוקיים .ii (בתרגול 11 ראינו את הדוגמה הזו אבל ללא סוגריים).

$$(x+y) * z \in L_2$$
$$x * (z * (z + y)) \in L_2$$

- : בהינתן דקדוקים שפת האיחוד שלהם נגדיר את  $G_1=\{V_1,\Sigma_1,R_1,S_1\},G_2=\{V_2,\Sigma_2,R_2,S_2\}$  נגדיר את בהינתן בהינתן  $L_3=L(G_1)\cup L(G_2)$  שימו לב שהתשובה צריכה להיות תלויה בפרמטרים של הדקדוקים  $G_1,G_2$
- בהן שבהן הבינאריות שבהן כל כל המחרוזות הבינאריות שבהן. כל באריות שבהן הבינאריות שבהן הבינאריות שבהן הבינאריות שבהן גער האפסים שווה למספר האחדות.

ב.

המקבלת generate\_language\_rec(rule\_dict, start\_var, k, mem) המקבלת השלימו את המימוש של הפונקציה (rule\_dict, start\_var, k, mem) ומשתנה ההתחלה CNF באותו האופן שבו CNF דקדוק G בצורת CNF המיוצג על ידי המילון set שמכיל את כל המילים באורך בשפה CYK. על הפונקציה להחזיר sets\_concat שמכיל את כל המימוש של הפונקציה הרקורסיבית. נתונה לכם פונקציית המעטפה, פונקציית העזר sets\_concat ומרבית המימוש של הפונקציה הרקורסיבית. הורידו את כל סימני ה-# מהקוד בקובץ השלד והשלימו את החסר היכן שמצוין (בכל מקום שצריך להשלים מתאימה שורת קוד אחת).

בסעיפים הבאים בחרו את התשובה ההדוקה ביותר. למשל, אם סיבוכיות הפונקציה היא לינארית וסימנתם שהיא לכל היותר פולינומיאלית, התשובה אינה הדוקה מספיק. ציינו בקובץ ה-PDF את התשובה שבחרתם **ונמקו בקצרה**. ניתן להניח שגודל האלפבית הוא קבוע.

- ii. סיבוכיות הזיכרון של הפונקציה generate\_language במקרה הגרוע היא:
  - .k-ביותר לינארית ב-(1
  - .k-ב לכל היותר פולינומיאלית ב-(2
  - .kכל הפחות אקספוננציאלית ב-(3
  - iii. סיבוכיות הזמן של הפונקציה generate language במקרה הגרוע היא:
    - .k- לכל היותר לינארית ב-(1
    - k-לכל היותר פולינומיאלית ב-(2
    - k-ב לכל הפחות אקספוננציאלית ב-(3

٦.

i. בקובץ השלד נתונה לכם הפונקציה (what(rule\_dict, start\_var, k הסבירו מילון ומספר טבעי k. הסבירו בקובץ השלד נתונה לכם הפונקציה (לא יותר משורה אחת) מה מחזירה הפונקציה.

בסעיפים הבאים בחרו את התשובה ההדוקה ביותר. למשל, אם סיבוכיות הפונקציה היא לינארית וסימנתם שהיא לכל היותר פולינומיאלית, התשובה אינה נכונה. ציינו בקובץ ה-PDF את התשובה שבחרתם **ונמקו בקצרה**. ניתן להניח שגודל האלפבית הוא קבוע.

- ii. סיבוכיות הזיכרון של הפונקציה what במקרה הגרוע היא:
  - .k-ב לכל היותר לינארית ב (1
  - k-ב לכל היותר פולינומיאלית ב-(2
  - k-ם לכל הפחות אקספוננציאלית ב-(3
  - iii. סיבוכיות הזמן של הפונקציה what במקרה הגרוע היא:
    - .k-ביותר לינארית ב-(1
    - .k- לכל היותר פולינומיאלית ב-(2
    - k-ב לכל הפחות אקספוננציאלית ב-(3

# שאלה 2 (15 נקי)

שאלה זו עוסקת בגנרטורים. בשאלה זו ניתן להתעלם משגיאות הקשורות לאי דיוק של floating point.

שמקבלת מספרים מטיפוס שמייצר אינסופי שמייצר פפt\_next\_sum(gen) א. ממשו את פונקציית הגנרטור פפt\_next\_sum(gen) א. ממשו את פונקציית האנרטור ( $\{a_n\}_{n=1}^\infty$ , ומחזירה גנרטור אשר בקריאת ה-next שונים מאפס קומחזירה גנרטור אשר בקריאת ה-st

$$\frac{1}{a_1} + \frac{1}{a_1 \cdot a_2} + \dots + \frac{1}{a_1 a_2 \cdot \dots \cdot a_n}$$

הניח הגרוע ולהשתמש ב-0(1) זיכרון. ניתן להניח הנחיה: על כל קריאת next לגנרטור לרוץ בזמן פון מקרה הגרוע ולהשתמש ב-0(1) זיכרון. ניתן להניח שפעולת next שפעולת אנרטור הקלט gen רץ בזמן ובמקום שפעולת

<u>נקודה למחשבה</u>: כבר נתקלתם בסכום מהצורה הזו <u>בתרגיל בית 3 שאלה 3אי,</u> בהקשר של חישוב נומרי. חשבו כיצד גנרטור כזה היה עוזר לכם במימוש ההוא (רשות – ממשו מחדש את הפונקציה מתרגיל בית 3 עם הגנרטור).

ב. ממשו את פונקציית הגנרטור ( $gen\_sequence$  כך שיתקיים שאם ניתן אותה כקלט לפונקציית הגנרטור מסעיף ב. e אי, קריאות e אי, קריאות מסעיף אי יתנו לנו קירוב הולך ומשתפר של המספר

$$0!=1$$
 במז:  $\sum_{i=0}^{\infty} rac{1}{i!} = e$  כאשר על פי הגדרה,

: דוגמת הרצה

## שאלה 3 (30 נקי)

שאלה זו עוסקת בעצי האפמן.

באורך מעל אייב בן בקורפוס תווים (משר כל תוn בקורפוס מופיע כתון גתון גתון מעל אייב בן אייב בן מעל אייב בקורפוס מחת.  $t\geq 2$ 

.corpus שנוצר הקורפוס של האפמן על הקורפוס H שנוצר כתוצאה מהרצת האלגוריתם של האפמן על הקורפוס

עבור כל אחד מארבעת הפריטים שלפניכם, ציינו בקובץ ה-PDF את ערכו כפונקציה של t,n, והסבירו בקצרה את תשובתכם. אם ישנו ערך יחיד, ציינו אותו במפורש. אם ישנו טווח ערכים אפשרי, תנו ערך תחתון וערך עליון את תשובתכם. אם ישנו לב שיש לתת תשובות מדויקות, ולא במונחי  $O(\cdot)$ .

<u>תזכורת</u>: גובה של עץ הוא אורך מסלול ארוך ביותר <u>בקשתות</u> מהשורש לעלה כלשהו בעץ. כמו כן משקל של עלה בעץ הוא שכיחות התו המתאים בקורפוס (היזכרו כיצד הוגדר בכיתה משקל של צומת פנימי בעץ).

- H מספר העלים בעץ.i
- ii. משקל השורש של העץ
  - iii. גובה העץ
  - iv מספר הצמתים בעץ.
- ב. נגדיר את ה**רמה ה-i** בעץ כקבוצת כל הצמתים במרחק בדיוק i מהשורש (כלומר הרמה ה-0 כוללת רק את השורש וכוי). השורש, הרמה ה-1 כוללת רק את בניו של השורש וכוי).

i-מו כן נגדיר את **משקל** הרמה הi להיות סכום המשקלים של הצמתים ברמה ה

הוכיחו או הפריכו את הטענות הבאות. אם הטענה נכונה, הסבירו בקצרה מדוע היא נכונה. אם הטענה אינה נכונה, ספקו דוגמה נגדית לקורפוס ולעץ שמתקבל ממנו באמצעות האלגוריתם שנלמד בכיתה, והסבר קצר מדוע הדוגמה מפריכה את הטענה.

- .i. לכל קורפוס corpus ולכל עץ האפמן H שנוצר מהקורפוס, לכל הרמות משקל זהה.
- ii. לכל קורפוס corpus ולכל עץ האפמן H שנוצר מהקורפוס, ישנתן שתי רמות שונות שמשקלן שונה.
- ג. בתרגול 12 ניתחנו את הסיבוכיות של הפונקציות של קידוד האפמן תחת ההנחה ש $|\Sigma|=0$ . בסעיף זה נרצה להבין מה משתנה בניתוח בכל אחד מהשלבים ללא הנחה זו.

נסמן ב- $\Sigma$  את האלפבית, ב-n את אורך הקורפוס corpus, ב-סמן את האלפבית, ב-b את אורך החודעה לסמן ב-b את אורך החודעה bits, בדומה לסימונים מהתרגול.

עבור כל אחת מהפונקציות הבאות שמופיעות במחברת תרגול 12, נתחו בקצרה את סיבוכיות **הזמן הממוצעת,** ללא ההנחה ש $|\Sigma|=0$ . ניתן להניח שפעולות אריתמטיות לוקחות זמן קבוע.

- char\_count(corpus) .i
- build huffman tree(char count dict) .ii
  - generate\_hcode(htree, prefix="") .iii
    - compress(text, hcode) .iv
    - decompress(bits, htree) .v

## שאלה 4 (25 נקי)

השאלה עוסקת בשינוי באלגוריתם למפל-זיו לדחיסת טקסט.

כזכור, הפונקציה LZW\_compress מחזירה את ייצוג הביניים של דחיסת למפל-זיו של המחרוזת text. למשל:

```
>>> LZW_compress("abcdabc")
['a', 'b', 'c', 'd', [4, 3]]
>>> LZW_compress("abab")
['a', 'b', 'a', 'b']
>>> LZW_compress("ababab")
['a', 'b', [2, 4]]

def inter_to_bin(intermediate, W=2**12-1, L=2**5-1)
```

שבהינתן רשימה lst עם ייצוג ביניים של מחרוזת דחוסה, מחזירה מחרוזת של ביטים, המייצגת את רצף הביטים לאחר הדחיסה. נזכיר, שתו שלא נדחס ייוצג ע"י הביט 0 ואחריו 7 ביטים עבור התו עצמו (סה"כ 8 ביטים, כלומר גם בתרגיל זה אנחנו מניחים לשם פשטות כי אנחנו מטפלים רק בתווי ASCII), ואילו מקטע שנדחס ייוצג ע"י הביט 1 ואחריו 12 ביטים עבור ההיסט אחורה, ו- 5 ביטים עבור אורך המקטע שנדחס (סה"כ 18 ביטים). שימו לב שבחישוב זה לקחנו בחשבון את ערכי ברירת המחדל של הפרמטרים W,L של שתי הפונקציות.

: דוגמאות הרצה

בעמוד הבא מופיעה הפונקציה שונה במעט עבור בעמוד הבא בעמוד הבא בעמוד במעט עבור במעט עבור בעמוד הבא מופיעה בעמוד במעט עבור בעמוד בעמוד בעמוד במעט עבור בעמוד בעמוד בעמוד במעט עבור בעמוד בעמוד בעמוד בעמוד בעמוד במעט עבור בעמוד בעמו

.LZW compress new, LZW compress ביחס ל: LZW compress ענו בקובץ ה $\operatorname{pd} f$  ענו בקובץ ה

- א. תנו דוגמא למחרוזת s המקיימת:
- .LZW\_compress(s) = LZW\_compress\_new(s)

מה יהיה הפלט (ייצוג הביניים) שיתקבל בשתי ההרצות?

ב. טענה: קיימת מחרוזת s שמקיימת:

```
len(inter to bin(LZW compress new(s))) < len(inter to bin(LZW compress(s)))</pre>
```

תנו דוגמא למחרוזת s כזו בצירוף שני ייצוגי הביניים המתקבלים עייי הפעלת כל אחת מהפונקציות הנייל, או הסבירו מדוע אין מחרוזת s כזו.

ג. טענה: קיימת מחרוזת s שמקיימת:

```
len(inter to bin(LZW compress new(s))) > len(inter to bin(LZW compress(s)))
```

תנו דוגמא למחרוזת s כזו בצירוף שני ייצוגי הביניים המתקבלים עייי הפעלת כל אחת מהפונקציות הנייל, או הסבירו מדוע אין מחרוזת s כזו.

```
def LZW compress new(text, start=0, W=2**12-1, L=2**5-1):
    n = len(text)
    if start >= n:
        return []
    #find the maximal length matching
    m,k = maxmatch(text, start, W, L)
    res1 = [text[start]] + \
         LZW compress new(text, start+1, W, L)
    res1 len = len(inter to bin(res1, W, L))
    if k < 3:
        return res1
    res2 = [[m,k]] + LZW compress new(text, start+k, W, L)
    res2 len = len(inter to bin(res2, W, L))
    if (res2 len < res1 len):
        return res2
    return res1
```

## שאלה 5 (20 נקי)

שאלה זו עוסקת בקודים לתיקון שגיאות.

א. לכל אחת מהטענות הבאות יש לסמן האם היא נכונה או לא. אם הטענה נכונה, יש להסביר בקצרה מדוע. אם היא לא נכונה, יש לספק דוגמא נגדית.

: הערות

- A-ם ב-רים מספר הוא מספר החיתוך שלהן ו-A הוא החיתוך האיברים ב-A הוא מספר האיברים ב-A
- C נאמר ש-C הוא קוד (n,k,d) אם  $(n,k,d)^k \to (0,1)^k$  היא פונקציה חחייע, ומרחק ההאמינג של  $\Delta \mathcal{C} = d$ 
  - d>1 קוד נקרא לא טריוויאלי אם •
  - על ידי y סביב r סביב את הכדור מגדירים את אנחנו עבור  $y \in \{0,1\}^n$  סביב על ידי כפי שראיתם בהרצאה, עבור  $y \in \{0,1\}^n$

כאשר  $\Delta$  הוא מרחק האמינג.

- $\left|B\left(y,\left\lfloor\frac{d-1}{2}\right\rfloor\right)\cap\mathrm{Im}\mathcal{C}
  ight|\leq 1$  מתקיים  $y\in\{0,1\}^n$  לא טריוויאלי אזי לכל (n,k,d) אם .i
  - $|B(y,\mathrm{d}-1)\cap\mathrm{Im}\mathcal{C}|=1$  מתקיים  $y\in Im\mathcal{C}$  אזי לכל טריוויאלי אזי לכל (n,k,d) אם .ii
    - $\ell > \left\lfloor \frac{d-1}{2} \right\rfloor$  אם  $\ell$  הוא קוד (n,k,d) לא טריוויאלי ו- $\ell$  מספר שלם המקיים כי .iii

 $|B(y,\ell)\cap \mathrm{Im}\mathcal{C}|>1$  מתקיים  $y\in\{0,1\}^n$  אזי לכל

- d=1 אזי בהכרח או (n,n-1,d) הוא הוא C אי .iv
  - d=1 אזי בהכרח או (n,n,d) או הוא C אם .v
- ב. בסעיף זה נתייחס לקודים מהצורה  $\mathcal{C}:\{0,1\}^k \to \{0,1,2,3\}^n$ , כלומר  $\mathcal{C}$  היא פונקציה חחייע שממירה מחרוזות בינאריות באורך  $\mathcal{C}:\{0,1\}^k \to \{0,1,2,3\}^n$  בינאריות באורך  $\mathcal{C}:\{0,1\}^k \to \{0,1,2,3\}^n$  באורך  $\mathcal{C}:\{0,1\}^k \to \{0,1,2,3\}^n$  בינאריות באורך למחרוזות בבסיס 4 באורך  $\mathcal{C}:\{0,1\}^k \to \{0,1,2,3\}^n$  בעבור כל אחד משני הקודים הבאים :
  - ערכל  $D\colon\{0,1,2,3\}^n\to\{0,1\}^k$  כלומר פונקציה ל-C, כלומר הפענוח ח המתאימה בענוח ח תארו את פונקציית הפענוח  $D\colon\{0,1,2,3\}^n\to\{0,1\}^k$  מתקיים השוויון  $x\in\{0,1\}^k$
- אם הקוד הוא טריוויאלי תנו דוגמה של שתי מילות קוד במרחק 1 זו מזו. אם הקוד אינו טריוויאלי, הוכיחו שמרחק ההאמינג של הקוד הוא לפחות 2.
- $b_i \in \{0,1\}$  מקבלת מחרוזות בינארית באורך  $p_i$ , נסמנה  $b_0 = c: \{0,1\}^n \to \{0,1,2,3\}^n$  .i נסמן ב- $p_i$  את המספר  $p_i = c: \{0,1\}$ , פורמלית  $p_i = c: \{0,1,2,3\}$  מחרוזת הקלט, פורמלית  $p_i = c: \{0,1\}$  כלומר  $p_i = c: \{0,1\}$  מוגדרת על ידי:  $p_i = c: \{0,1\}$  מוגדרת על ידי:

$$C(b_0b_1\dots b_{n-1})=tb_2\dots b_{n-1}\textcolor{red}{p}$$

n=6 לדוגמה עבור

$$C(110010) = 300101$$
  
 $C(011000) = 110000$ 

n=6 לדוגמה עבור

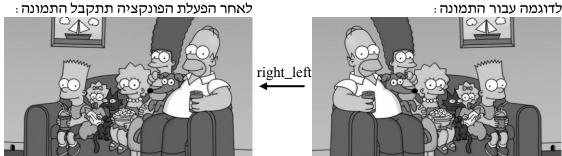
$$C(110010) = 300100$$
  
 $C(011000) = 110002$ 

## <u>שאלה 6 (10 נקי)</u>

שימו לב שעל מנת להריץ את הפונקציות בשני הסעיפים הבאים עליכם להתקין את הספרייה PILLOW על ידי הרצת הפקודות הבאות ב-python3 (אם זה לא עובד, נסו להחליף את המילה python3 ב-python3. אם אתם עדיין נתקלים בבעיות, היעזרו בפיאצה ובשעות החונכות):

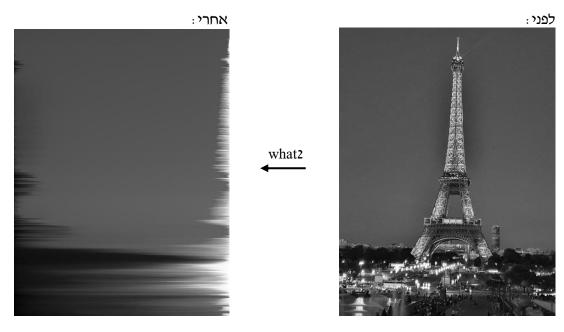
python -m pip install --upgrade pip python -m pip install --upgrade Pillow

א. השלימו בקובץ השלד את שלוש השורות החסרות בפונקציה right\_left, שמקבלת מטריצה שמייצגת תמונה ומחזירה מטריצה חדשה שמייצגת את התמונה שמתקבלת ע"י שיקוף התמונה על הציר האנכי.



בדקו את הפונקציה שלכם על ידי הפעלתה על התמונה the-simpsons.jpg המצורפת בין קבצי התרגיל.

ב. לפניכם תמונה לפני ואחרי שהפעלנו עליה פעולה כלשהי.



השלימו בקובץ השלד את הפונקציה what2 על ידי הוספת שתי שורות מתאימות כדי להגיע לאותה תוצאה. בדקו את הפונקציה שלכם על ידי הפעלתה על התמונה eiffel-tower.jpg המצורפת בין קבצי התרגיל.