

מ.מ.למדמ"ח ~ משהו?

שחר פרץ

29 במאי 2024

1 דקדוק - תזכורות

תזכורת:

$$G = (V, \Sigma, R, S)$$

- V - משתנים
- S_i - אלפבית
- R - כללי גזירה
- $S \in V$ - משתנה התחלה

דוג':

$$V = \{A, B\} \quad (1)$$

$$\Sigma = \{0, 1, \#\} \quad (2)$$

$$R: A \rightarrow 0A1 \iff A \rightarrow 0A1 \mid B \quad (3)$$

$$A \rightarrow B \quad B \rightarrow B$$

$$B \rightarrow B$$

$$S = A \quad (4)$$

במקרה הזה:

$$L(G) = \{0^n \# 1^n \mid n \in \mathbb{N}\}$$

לדוגמה $00\#11 \in L(G)$. תזכורת: אנו נתמקד אך ורק בדקדוקים חסרי הקשר - בכללי הגזירה, בצד שמאל מפיע רק משתנה אחד (וללא בתווים מ- Σ). לדוגמה, דקדוק שיכלול את הכלל $aV \rightarrow A$, יהיה לא חסר הקשר. (חסר הקשר = Context Sensitive ודקדוק לא חסר הקשר = Context Free Grammars).

שאלות ששאלנו בשיעור שעבר:

$$\bullet \text{ האם } L(G_1) = L(G_2)?$$

$$\bullet \text{ האם } L(G_1) \subseteq L(G_2)?$$

$$\bullet \text{ האם יש יותר מעץ גזירה אחד לכל ילה ב- } L(G_1)?$$

$$\bullet \text{ האם } L(G_1) \cap L(G_2) \neq \emptyset?$$

כל הדברים להלן, שקולים לבקשה לכתוב תוכנת מחשב המקבלת כקלט תוכנית מחשב, ולהחזיר אם התוכנית תיגמר בזמן סופי או לא. השאלה הזו בלתי אפשרית לפתרון סופי.

השאלות האלו קרויות **לא כריעות** - מלשון הכרעה, לא ניתן להכריע את התוצאה שלהן.

נרצה לענות בזמן טוב על שאלות יותר קלות: בהינתן דקדוק חסר הקשר G , ומילה w :

$$\bullet \text{ האם } w \in L(G) \text{ (Recognition)?}$$

$$\bullet \text{ אם } w \in L(G), \text{ למצוא עץ גזירה כלשהו של } w \text{ ב- } G \text{ (Parsing).}$$

2 פתרון הבעיות

2.1 CNF

על מנת לענות על השאלות הללו, נמיר דקדוק ל-Chomsky Normal Form, ובצורה זו קל יותר לענות על השאלות האלו.
הגדרה: דקדוק חסר הקשר הוא מצורת CNF אם"מ כל כללי הגזירה בו הם מהצורה הבאה:

$$A \rightarrow a, \quad A \in V, a \in \Sigma \quad (5)$$

$$A \rightarrow BC \quad A \in V, B, C \in V \setminus \{S\} \quad (6)$$

$$S \rightarrow \varepsilon \quad (7)$$

שימו לב שזו הגדרה הנכונה, בזו של השיעור שעבר נפלה טעות

טענה: כל דקדוק חסר הקשר ניתן לכתיבה והמרה בצורת CNF.

כדי "להוכיח", נראה כיצד נפתור כללי גזירה לא תקינים:

$$A \rightarrow BCD \implies \begin{cases} A \rightarrow BE \\ E \rightarrow CD \end{cases} \quad (8)$$

$$A \rightarrow \varepsilon \implies \begin{cases} B \rightarrow A \\ B \rightarrow a|aA \end{cases} \quad (9)$$

$$A \rightarrow aB \implies \begin{cases} A \rightarrow A'B \\ A' \rightarrow a \end{cases} \quad (10)$$

נתרגם את בחלק (1) למעלה ל-CNF. שני הכללים הראשונים לא עונים להגדרה.

$$\begin{cases} A \rightarrow 0A1 \\ A \rightarrow B \\ B \rightarrow \# \end{cases} \implies \begin{cases} C_0 \rightarrow 0 \\ C_1 \rightarrow 1 \\ A \rightarrow C_0AC_1 \\ A \rightarrow B \\ B \rightarrow \# \end{cases} \implies \begin{cases} C_0, C_1 \rightarrow 0, 1 \\ A \rightarrow C_0A \\ A_1 \rightarrow AC_1 \\ A \rightarrow \# \end{cases} \quad (11)$$

2.1.1 סיכום

בצורת CNF עלולים להיות יותר משתנים ויותר כללי גזירה, ובכלל, הדקדוק עלול להיות פחות קריא לאנשים. אבל, עבור המחשב, צורה זו מאוד נוחה.

2.2 אלגוריתם CYK – Cocke-Younger-Karzem

אלג' לזיהוי ופירוש של מילה בהינתן דקדוק חסר הקשר בצורת CNF.

רעיון (רקורסיבי): אם $G = (V, \Sigma, R, S)$ ונתונה המילה w , נאמר $w \in L(G)$ אם"מ מתקיימים אחד משלושת התנאים הבאים:

- $w = \varepsilon$ וקיים הכלל $S \rightarrow \varepsilon$.
- $w = x \in \Sigma$ וקיים הכלל $S \rightarrow x$.
- קיימת חלוקה $1 \leq k \leq |w| - 1$ של w לשתי מחרוזות לא ריקות $w[1:k]$, $w[k+1:|w|]$ וקיים כלל $S \rightarrow XY$ כך שניתן לגזור את $w[1:k]$ מ- X וניתן לגזור את $w[k+1:|w|]$ מ- Y .

זהו לא פתרון יעיל, אך הוא עובד, וסופי, ונוכל ליעל אותו.

לצורך נוחות, נסמן מתשנים באותיות גדולות באנגלית ותווים בעזרת אותיות קטנות, ואת כללי הגזירה נייצג במילון עם משתנה לקבוצת הכללים שלו. לדוגמה, $A \rightarrow AB$, $A \rightarrow a$, $B \rightarrow b$, יהיה שקול ל-:

$$R = \{A: \{AB, a\}, B: \{b\}\}, S = \{A, B\}, V = \{a, b\}, S = A$$

זה הזמן להתחיל לחרבש קוד:

```

1 def CYK(rules, start_var, w):
2     if len(w) < 2:
3         return w in rules[start_var]
4     for k in range(1, len(w)):
5         for rule in rules[start_var]:
6             if len(rule) == 2:
7                 if CYK(rules, rule[0], w[:k]) and \
8                     CYK(rules, rule[1], w[k:]):
9                     return True
10
11 return False

```

דוג':

$$V = \{S, A < B < C\} \quad (12)$$

$$\Sigma = \{a, b\} \quad (13)$$

$$R: S \rightarrow AB|BC \quad (14)$$

$$A \rightarrow BA|a$$

$$B \rightarrow CC|b$$

$$C \rightarrow AB|a$$

$$w = baaba \quad (15)$$

נמצא שכבר בקריאה הראשונה לרקורסיה, נעשה 16 קריאות רקורסיביות עוד בקריאה הראשונה לפונקציה. לא אפרט איך זה עובד למאות שהמורה פירט את זה על הלוח כי זה נראה לי מיותר.

יש כאן פעולה מאוד מיותרת ויקרה - slicing. נוכל לשנות את חתימת הפונקציה, כך שהיא תקבל אינדקס של i, j אינטרוואל קריאה. כדי לחסוך את זה - ראה הערות על הקוד.

הסיבוכיות של המימוש הרקורסיבי היא לפחות אקספוננציאלית. הסיבה - אנחנו מחשבים שוב ושוב את אותם הערכים. נרצה להבדיל בין תכנון דינמי לממואיזה - בממואיזה נביא פטיש, נשמור הכל מחוץ לרקורסיה ונקווה לטוב, ובתכנון דינמי נתכנן את הפתרון ונבנה את זה באמצעות ממואיזה.

כמות עצי הגזירה השונים - אספוננציאלית (הרי זה שקול ללעבור באלג' הרגיל בלי הממואיזה). אך כמות הצמתים בעץ היא כמעט n . נשתמש בעובדה שעשינו ממואיזה (פשוט תניחו שהקוד כתוב איפשהו). נגדיר $rr = |rules|$ ו- $n = |w|$ אז משום שאנו צריכים לבחור בזכרון את $i, j, start_var$, כאשר $i \in [n], j \in [n]$ ו- r חוסם את כמות האפשרויות של $start_var$, אז סה"כ גודל הטבלה יהי $n^2 r$ לכל היותר. זמן הריצה יהיה:

$$O\left(\sum_{i < j} \sum_{var} \sum_{i < k < j} |rules[var]|\right) = O(n^3 r^2)$$

אך זה לא הדוק. נוכל לשנות את הסדר של הסיגמות (זה טריק של סכומים שכדאי להכיר):

$$O\left(\sum_{i < j} \sum_{i < k < j} \sum_{var} |rules[var]|\right) = O\left(\sum_{i < k < j} \underbrace{\sum_{var} |rules[var]|}_{O(r)}\right) = O(n^3 r)$$

אומנם $rules[var] \leq r$ וגם כמות האפשרויות ל- var קטנה מ- r , ולכן זה חסום מלמעלה גם ע"י r^2 . אך למעשה בשני הסכומים הללו עברנו "דק" על כל החוקים (כל var וכל אחד מהחוקים שלו)!