

Algorithme de Cocke-Younger-Kasami

Référence : 131 Développements pour l'oral, D. Lesesvre

Introduction. L'algorithme CYK permet de résoudre le problème du mot pour les grammaires algébriques. On dit qu'une grammaire $G = (V, \Sigma, R, S)$ est sous forme normale de Chomsky si elle ne contient que des règles de la forme $A \rightarrow a$ (avec $a \in \Sigma$ et $A \in V$), ou $A \rightarrow A_1 A_2$ ($A_1, A_2 \in V$).

Définition 1 (Problème de mot). *Étant donné un mot $w = w_1 \dots w_n$ et une grammaire G sous forme normale de Chomsky, a-t-on $w \in L(G)$?*

Algorithme CYK. On utilisera de la programmation dynamique. Pour tous les indices $1 \leq i \leq j \leq n$, on note

$$E_{i,j} = \{A \in V, A \rightarrow^* w_i \dots w_j\}$$

Ainsi, $w \in L(G)$ ssi $S \in E_{1,n}$. On va donc calculer $E_{1,n}$ par programmation dynamique.

Initialisation. Montrons que pour tout $1 \leq i \leq n$, on a

$$E_{i,i} = \{A \in V : A \rightarrow w_i \in R\}$$

En effet, si $A \in E_{i,i}$, alors $A \rightarrow^* w_i$. Or, puisque G est sous forme normale de Chomsky, on a jamais $X \rightarrow^* \epsilon$ pour $X \in V$. Ainsi, la première règle de $A \rightarrow^* w_i$ n'est pas de la forme $A \rightarrow A_1 A_2$ qui donnerait un mot d'au moins 2 lettres. On en déduit qu'on applique une règle de la forme $A \rightarrow a$, et donc $a = w_i$. Ainsi, $A \rightarrow w_i \in R$.

L'autre inclusion est triviale.

Récurrence. Soit $1 \leq i < j \leq n$, montrons

$$E_{i,j} = \bigcup_{k=i}^{j-1} \bigcup_{\substack{B \in E_{i,k} \\ C \in E_{k+1,j}}} \{A \in V : A \rightarrow BC \in R\}$$

Montrons l'inclusion indirecte. Soit $A \in V$ tel qu'il existe k tel que $i \leq k \leq j-1$, $B \in E_{i,k}$ et $C \in E_{k+1,j}$ avec $A \rightarrow BC \in R$.

On a alors $B \rightarrow^* w_i \dots w_k$, et $C \rightarrow^* w_{k+1} \dots w_j$. Ainsi, on a :

$$A \rightarrow BC \rightarrow^* w_i \dots w_k C \rightarrow^* w_i \dots w_j$$

Ainsi, $A \in E_{i,j}$.

Réciproquement, soit $A \in E_{i,j}$, c'est-à-dire $A \rightarrow^* w_i \dots w_j$. On a au moins deux lettres dans le mot $w_i \dots w_j$, la première règle appliquée est donc de la forme $A \rightarrow BC$. On utilise alors l'arbre de dérivation de $A \rightarrow^* w_i \dots w_j$, A est la racine et a deux enfants, B et C . Pour exhiber k , il suffit de prendre le nombre l de feuilles dans l'arbre enraciné en B , et on pose $k = i + l - 1$. Ainsi, on a $B \rightarrow^* w_i \dots w_k$ et $C \rightarrow^* w_{k+1} \dots w_j$. Cela conclut la preuve.

Algorithme. L'algorithme va donc d'abord calculer les $A_{i,i}$ et remonter jusqu'à l'ensemble $E_{1,n}$ en calculant diagonale par diagonale.

Algorithme 1 : CYK(w, G)

```

pour  $1 \leq i \leq j \leq n$  faire
     $E_{i,j} \leftarrow \emptyset$ ;
pour  $i = 1 \dots n$  faire
    pour  $A \rightarrow a \in R$  faire
        si  $a = w_i$  alors
             $A_{i,i} \leftarrow E_{i,i} \cup \{A\}$ ;
pour  $d = 2 \dots n$  faire
    pour  $(i, j)$  sur la d-diagonale supérieure faire
        pour  $k = i \dots j - 1$  faire
            pour  $A \rightarrow B_1 B_2$  faire
                si  $B_1 \in E_{i,k}$  et  $B_2 \in E_{k+1,j}$  alors
                     $E_{i,j} \leftarrow E_{i,j} \cup \{A\}$ ;
retourner  $S \in E_{1,n}$ ;
    
```

Implémentation et complexité. On veut pouvoir ajouter et vérifier rapidement si des non-terminaux appartiennent à un ensemble $E_{i,j}$. On peut alors, pour tout $1 \leq i \leq j \leq n$, utiliser un tableau booléen de taille $|V|$ pour avoir les deux opérations en temps constant.

On a alors une complexité temporelle en $\mathcal{O}(|R| \times n^3)$. pour la complexité spatiale, on a besoin de $\frac{n(n+1)}{2}$ tableaux de taille $|V|$, donc un $\mathcal{O}(n^2|V|)$.

Commentaires. Il faut d'abord transformer une grammaire pour la mettre sous forme normale de Chomsky. Cette transformation peut faire exploser la taille de la grammaire. Il faut aussi vérifier qu'on a bien $\epsilon \notin L(G)$, ce qui peut se faire par saturation.