

CYK Algorithm

Shamshad Ali

Computer Science & Engineering

Babu Banarasi Das Institute of Engineering Technology & Research Centre, Bulandshahr, U.P

shamshadali02@gmail.com

ABSTRACT

In this paper, we have to test the string that the string is a part of the given language. We have the substrings but we are not conformed that all these belong to given language. So how can we say that this string is that part of the given language or not. Now we are confused to determine it. To solve this problem we have a lot of way to check the membership of the given string. So that, we come to know about the membership of the sub-string or a string.

So to check the membership we have the algorithm that was invented by John Cocke, Taddo Kasami (in 1965) and Daniel H. Younger (in 1967). Although from the starting it was not so much popular. But now it is in used. It is time consuming. Because it takes too much time to solve and check the string whether it is the member or not. The reason behind this is the equation. Every step takes the help of equation. It is step by step process. If one step is calculated wrong then all the next will be wrong. So it takes more attention while computing.

It clearl-y checks the membership by the help of equation that value is given in the table. So we will test by the help of the table. If the starting symbol seems in the table entry V_{1n} , then the string is the member of the language otherwise no. The whole table construction process takes $O(n^3)$ time in worst case. If w is of length n then there will be exactly $2n-1$, nodes labelled by variables in the tree.

Keywords: Cocke-Younger-Kasami (CYK), Context free language (CFL), Chomsky Normal Form (CNF), Grammar (G).

I. INTRODUCTION

Basically the CYK algorithm is used to check or test the string whether the string belongs to the given language or not(i.e. the given string, is the member of the given language). We can describe membership of a string w in a **CFL**¹ L . There is an efficient technique based on the idea of “Dynamic Programming” which may known as “Table Filling Algorithm” or

“Tabulation”. This algorithm known as CYK Algorithm (i.e. **Cocke-Younger-Kasami**).

The algorithm works only if the grammar is in Chomsky normal form (CNF)² and succeeds by breaking one problem into a sequence of smaller one. Assume that we have a CNF grammar

$G = (V, T, P, S)$ ³ for a CFL language L and a string

$$w = a_1 a_2 a_3 \dots\dots\dots a_n \quad \text{in } T^*$$

and the substrings are defined as

$$w_{ij} = a_i \dots\dots\dots a_j$$

and the subsets that will be used in the table are defined as

$$V_{ij} = \{ A \in V : A \Rightarrow^* w_{ij} \}$$

It is simple that, $w \in L(G)$ if and only if $S \in V_{1n}$. The algorithm constructs a table that tells whether w is in L . Note that when we will compute this w , the grammar itself is considered fixed, and its size contributes only a constant factor at the running time, which is measured in terms of the length of the string w whose membership in L is being tested.

In the CYK algorithm, we construct a triangular table as shown in the below figure 1.1

* CFL, CNF and G are defined at the last before the abbreviation used page.

Row 5	V_{15}				
Row 4	V_{14}	V_{25}			
Row 3	V_1	V_{24}	V_{35}		
Row 2	V_{12}	V_{23}	V_{34}	V_{45}	
Row 1	V_{11}	V_{22}	V_3	V_{44}	V_{55}
	a_1	a_2	a_3	a_4	a_5

Figure 1: CYK table.

The horizontal axis corresponds to the positions of the string, $w = a_1 a_2 a_3 \dots\dots\dots a_n$, which we have

supposed has length 5. The table entry (or Elements of the Table) V_{ij} is the set of variables A such that:

$$A \Rightarrow a_i a_{i+1} \dots a_j$$

To compute V_{ij} , observe that $A \in V_{ij}$ if and only if G contains a production $A \rightarrow a_i$. Therefore V_{ij} can be computed for all $(1 \leq i \leq n)$ by inspection of w and the production of the grammar. To continue, notice that for $j > i$, A derives w_{ij} if and only if there is a production $A \rightarrow BC$, with $B \Rightarrow w_{ik}$ and $C \Rightarrow w_{k+1,j}$ for some k with $i \leq k$, and $k < j$. In other words $V_{ij} = \cup \{ A : A \rightarrow BC, \text{ with } B \in V_{ik} \text{ and } C \in V_{k+1,j} \mid k \in \{i, i+1, \dots, j-1\} \}$ Equation 1.1

Row 1, Row 2, Row 3 are also computed as
 Row 1: compute $V_{11}, V_{22}, \dots, V_{nn}$
 Row 2: Compute $V_{12}, V_{23}, \dots, V_{n-1,n}$
 Row 3: Compute $V_{13}, V_{24}, \dots, V_{n-2,n}$
 and so on.

To fill the table, we work row by row, upwards. Notice that each row corresponds to one length of substrings; the bottom row is for strings of length 1, the second from bottom row for strings of length 2, and so on until the top row corresponds to the one substring of length n , which is w itself. It takes $O(n)$ time to compute any one entry of the table by a method. Since there are $n(n+1)/2$ table entries, the whole table construction process takes $O(n^3)$ time in worst case. As we know Chomsky normal form (CNF) grammar are binary trees, if w is of length n then there will be exactly $2n-1$ nodes labeled by variables in the tree. Example 1.1 The CNF grammar is defined a

$$\begin{aligned} S &\rightarrow AB / BC & A &\rightarrow BA / a \\ B &\rightarrow CC / b & C &\rightarrow AB / a \end{aligned}$$

Suppose $w = baaba$ i.e. of length $n = 5$. If we will draw the binary tree, then there will be exactly $2n-1 = 2*5 - 1 = 9$ nodes labeled variables. Let us see by constructing tree

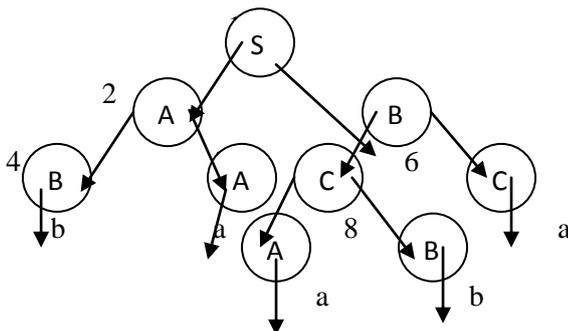


Figure 2: Binary tree

Here $w = "baaba"$ i.e. length $n=5$. Then no. of nodes must be $2n-1 = 2*5-1=10-1=9$. As shown in the figure there are 9 nodes variables. In any condition if we draw binary tree either left most or right most of fixed length then there nodes always be as according to $2n-1$.

Note: If we draw the table for this string $w = baaba$ of length 5 then there will be total $n(n+1)/2$ entries will be found.

$$\text{So, } n(n+1)/2 = 5(5+1)/2 = 15 \text{ entries.}$$

V_{15}				
V_{14}	V_{25}			
V_{13}	V_{24}	V_{35}		
V_{12}	V_{23}	V_{34}	V_{45}	
V_{11}	V_{22}	V_{33}	V_{44}	V_{55}

Figure 3: CYK Simple Table

Total 15 entries and the values V_{ij} 's are computed by the equation 1.1.

II. ALGORITHM

Let the input be a string S consisting of n characters: $a_1 \dots a_n$.
 Let the grammar contain r nonterminal symbols $R_1 \dots R_r$.
 This grammar contains the subset R_s which is the set of start symbols.
 Let $P[n,n,r]$ be an array of Booleans. Initialize all elements of P to false.

```

for each i=1 to n
    for each unit production  $R_j \rightarrow a_i$ 
        set  $P[i,1,j] = \text{true}$ 
    for each i=2 to n .....length of span
        for each j=1 to n-i+1 .....Start of span
            for each k=1 to i-1 .....partition of span
                for each production  $RA \rightarrow RB RC$ 
                    if  $P[j,k,B]$  and  $P[j+k, i-k, C]$ 
                        then set  $P[j,i,A] = \text{true}$ 
    if any of  $P[1,n,x]$  is true (x is iterated over the set s, where s are all the indices for Rs)
        then S is member of language
    else
        S is not member of language.
    
```

Question 1: Use the CYK algorithm to determine whether the string $w = baaba$ is in the language generated by the CNF grammar:

$$S \rightarrow AB / BC$$

$$\begin{aligned} A &\rightarrow BA / a \\ B &\rightarrow CC / b \\ C &\rightarrow AB / a \end{aligned}$$

Solution:

STEP 1: We have to compute all the values of V_{ij} 's .
That is Row 1 such as;

V_{11} is 'b' and 'b' is produced by $B \rightarrow b$. So $V_{11} = \{ B \}$, i.e. take the left side variable.

V_{22} is 'a' and 'a' is produced by $A \rightarrow a$ and $C \rightarrow a$ so we take union of A and C. so

$$V_{22} = \{ A, C \}. \text{Similarly}$$

$$V_{11} = V_{44} = \{ B \} \text{ and } V_{22} = V_{33} = V_{55} = \{ A, C \}$$

STEP 2: Now we will go upward that is we will compute Row2. We have to take the help of equation 1.1 ($A \rightarrow BC$).

So, $V_{12} \rightarrow V_{ik} V_{k+1,j}$ here $k=1$ only so $V_{ik} \in V_{11}$ and $V_{k+1,j} \in V_{22}$. Put these values we find;

$$V_{12} \rightarrow V_{11} V_{22}. \text{ Put the values of } V_{11} \text{ and } V_{22}.$$

$$\rightarrow \{ B \} \{ A, C \}$$

$\rightarrow \{ BA, BC \}$ BA and BC are produced by $S \rightarrow BC$ and $A \rightarrow BA$

$$V_{12} \rightarrow \{ S, A \}.$$

Similarly,

$V_{23} \rightarrow V_{ik} V_{k+1,j}$ here $k=2$ only so $V_{ik} \in V_{22}$ and $V_{k+1,j} \in V_{33}$

$$\rightarrow V_{22} V_{33}$$

$$\rightarrow \{ A, C \} \{ A, C \}$$

$$\rightarrow \{ AA, AC, CA, CC \}$$

AA, AC, CA, are not in the grammar.

But CC is in the grammar and produced by $B \rightarrow CC$

$$\text{So, } V_{23} \rightarrow \{ B \}.$$

Similarly,

$$V_{34} \rightarrow V_{ik} V_{k+1,j} \text{ here } k=3$$

so $V_{ik} \in V_{33}$ and $V_{k+1,j} \in V_{44}$

$$\rightarrow V_{33} V_{44}$$

$$\rightarrow \{ A, C \} \{ B \}$$

$$\rightarrow \{ AB, CB \}$$

CB is not in the grammar.

AB is in the grammar and produced by $S \rightarrow AB$ and $C \rightarrow AB$

$$\text{So, } V_{34} \rightarrow \{ S, C \}.$$

Similarly,

$$V_{45} \rightarrow V_{44} V_{55} \text{ here } k=4$$

$$\rightarrow \{ B \} \{ A, C \}$$

$$\rightarrow \{ BA, BC \}$$

BA is produced by $A \rightarrow BA$ and BC is produced by $S \rightarrow BC$

$$V_{45} \rightarrow \{ S, A \}.$$

STEP 3: Now Row 3 will be computed

$$V_{13} \rightarrow V_{ik} V_{k+1,j} \text{ here } k=1,2$$

$$\text{So for } k=1 \quad V_{13} \rightarrow V_{11} V_{23}$$

$$\text{for } k=2 \quad V_{13} \rightarrow V_{12} V_{33}$$

So we will take union of both. i.e.

$$V_{13} \rightarrow V_{11} V_{23} \cup V_{12} V_{33}$$

$$\rightarrow \{ B \} \{ B \} \cup \{ S, A \} \{ A, C \}$$

$\rightarrow \{ BB, SA, SC, AA, AC \}$ Nothing are in the grammar so it is null.

$$V_{13} \rightarrow \{ \emptyset \}$$

Similarly,

$$V_{24} \rightarrow V_{22} V_{34} \cup V_{23} V_{44}$$

because here $k=2$ and 3 .

$$\rightarrow \{ A, C \} \{ S, C \} \cup \{ B \} \{ B \}$$

$$\rightarrow \{ AS, AC, CS, CC, BB \}$$

here only CC is produced by $B \rightarrow CC$

$$\text{So, } V_{24} \rightarrow \{ B \}$$

Similarly,

$$V_{35} \rightarrow V_{33} V_{45} \cup V_{34} V_{55}$$

here $k=3$ and 4 .

$$\rightarrow \{ A, C \} \{ S, A \} \cup \{ S, C \} \{ A, C \}$$

$$\rightarrow \{ AS, AA, CS, CA, SA, SC, CA, CC \}$$

here only CC is produced by $B \rightarrow CC$

$$\text{So, } V_{35} \rightarrow \{ B \}.$$

STEP 4: Now we will compute Row 4 as above.

$$V_{14} \rightarrow V_{ik} V_{k+1,j} \text{ here } k \text{ contains three values}$$

i.e. $k=1,2$ and 3 .

$$\text{So for } k=1 \quad V_{14} \rightarrow V_{11} V_{24}$$

$$\text{For } k=2 \quad V_{14} \rightarrow V_{12} V_{34}$$

$$\text{For } k=3 \quad V_{14} \rightarrow V_{13} V_{44}$$

$$V_{14} \rightarrow V_{11} V_{24} \cup V_{12} V_{34} \cup V_{13} V_{44}$$

$$\rightarrow \{ B \} \{ B \} \cup \{ S, A \} \{ S, C \} \cup \{ \emptyset \} \{ B \}$$

$$\rightarrow \{ BB, SS, SC, AS, AC, \emptyset \}$$

Nothing are produced.

$$\text{So, } V_{14} \rightarrow \{ \emptyset \}.$$

Similarly,

$$V_{25} \rightarrow V_{22} V_{35} \cup V_{23} V_{45} \cup V_{24} V_{55}$$

$$\rightarrow \{ A, C \} \{ B \} \cup \{ B \} \{ S, A \} \cup$$

$$\{ B \} \{ A, C \}$$

$$\rightarrow \{ AB, CB, BS, BA, BC \}$$

Here AB is produced by $S \rightarrow AB$ and $C \rightarrow AB$

BA is produced by $A \rightarrow BA$

BC is produced by $S \rightarrow BC$

So after taking union of left context i.e. A, S, and C.

$$V_{25} \rightarrow \{ S, A, C \}.$$

STEP 5: Now the last Row 5.

$$V_{15} \rightarrow V_{ik} V_{k+1,j} \text{ here the value of } k \text{ is five}$$

i.e. $k=1,2,3$ and 4 .

$$\text{So for } k=1 \quad V_{15} \rightarrow V_{11} V_{25}$$

$$\text{for } k=2 \quad V_{15} \rightarrow V_{12} V_{35}$$

$$\text{for } k=3 \quad V_{15} \rightarrow V_{13} V_{45}$$

$$\text{for } k=4 \quad V_{15} \rightarrow V_{14} V_{55}$$

$$V_{15} \rightarrow V_{11} V_{25} \cup V_{12} V_{35} \cup V_{13} V_{45} \cup V_{14} V_{55}$$

$$V_{15}$$

$$\rightarrow \{ B \} \{ S, A, C \} \cup \{ S, A \} \{ B \} \cup \emptyset \cup \emptyset$$

→ {BS, BA, BC, SB, AB} BS and SB are not in the grammar.
 BC is produced by $S \rightarrow BC$
 AB is produced by $S \rightarrow AB$ and $C \rightarrow AB$
 BA is produced by $A \rightarrow BA$

So, $V_{15} \rightarrow \{S, A, C\}$. Now put all the values in the table.

{S,A,C}					
{φ}	{S,A,C}				
{φ}	{B}	{B}			
{S,A}	{B}	{S,C}		{S,A}	
{B}	{A,C}	{A,C}	{B}	{A,C}	

Figure 4 Final CYK Table.

As we have got that the value of V_{15} is {S,A,C}. It contains 'S' which is the starting symbol. i.e. $S \in V_{15}$. So $w = "baaba"$ is the member of the given grammar.

Exercise: Check out the membership for the string $w = "aabb"$ by the CYK algo.
 $S \rightarrow AB$ $A \rightarrow BB/a$ and $B \rightarrow AB/a$.

CFL¹: It is also called as Type-2 Production. The production is of the form $\alpha \rightarrow \beta$ where $\alpha \in V$ and $\beta \in (V \cup \Sigma)^*$. For example, $S \rightarrow aSb$, $A \rightarrow a/b$, $A \rightarrow \epsilon$ etc. It is also called a context-free grammar. A language generated by a context-free grammar is called a context free language (CFL).

CNF²: In this form we have restrictions on the length of R.H.S. i.e. a Context free grammar G, is in CNF if every production is of the form

- 1.) $A \rightarrow BC$, where $A, B, C \in V$.
i.e. One non-terminal \rightarrow one non-terminal. One non-terminal
 - 2.) $A \rightarrow a$ where $A \in V$ and $a \in T$
One non-terminal \rightarrow Terminal
- Then this normal form is called CNF.

OR

For a grammar in CNF, the derivation tree has the following property: Every node has at most two descendents - either two internal vertices or a single leaf.

G³: Noam Chomsky gave the definition of a grammar. It is a simple 4-tuple language i.e. $G = (V, \Sigma, P, S)$, where

- 1.) V is a finite nonempty set whose elements are called Variables or Non-terminals. Mostly Capital letters are used. i.e. $V = \{A, B, X, Y, \dots\}$.

- 2.) Σ is a finite nonempty set whose elements are called Terminals. Mostly small letters are used. It is also denoted by 'T'.

So Σ or $T = \{a, b, c, \dots, x, y, \dots\}$.

- 3.) $V \cap \Sigma = \{\phi\}$.

- 4.) S is a special variable (an element of V i.e. $S \in V$), or called start symbol.

- 5.) P is a finite set. The elements of p are called Productions.

i.e. $P : S \rightarrow aS / bS / a/b$.

III. ABBREVIATION USED

- CYK : Cocke-Younger-Kasami
- CFL : Context-free language
- CFG : Context -free grammar
- CNF : Chomsky normal form
- G : Grammar
- V : Variable
- P : Production
- S : Special Symbol
- T or Σ : Terminal
- V_{ij} : Table Entry
- W : String
- W_{ij} : Substring
- K : Integer
- I : Integer
- J : Integer
- A,S,B,.. All capital letters are Variables or non-terminals
- a,b,c,... All small letters are Terminals.
- α : Left context
- β : Right context
- n : Natural no.
- $O(n^3)$: Worst case running time.
- {φ} : Null
- ϵ : Belongs to

REFERENCES

1. Hopcroft, Ullman, "Introduction to Automata Theory, Languages and Computation", Pearson Education.
2. K.L.P. Mishra and N. Chandrasekaran, "Theory of Computer Science: Automata, Languages and Computation", PHI
3. Peter Linz, "An Introduction to Formal Language and Automata", Narosa.
4. <http://nitishkr.wordpress.com/2011/03/29/cyk-algorithm.implementation>