

Lecture 15: CYK Parsing Algorithm

3 March 2009

1 CYK parsing

1.1 Discussion

We already saw that one can decide if a word is in the language defined by a context-free grammar, but the construction made no attempt to be practical. There were two problems with this algorithm. First, we converted the grammar to be in Chomsky Normal Form (CNF). Most practical applications need the parse to show the structure using the original input grammar. Second, our method blindly generated all parse trees of the right size, without regard to what was in the string to be parsed. This is inefficient since there may be a large number of parse trees (i.e., exponential in the length of the word) of this size.

The Cocke-Younger-Kasami (CYK) algorithm solves the second of these problems, using a table data-structure called the *chart*. This basic technique can be extended (e.g. Earley's algorithm) to handle grammars that are not in Chomsky Normal Form and to linear-time parsing for special types of CFGs.

In general, the number of parses for a string w is exponential in the length of w . For example, consider the sentence “Mr Plum kill Ms Marple at midnight in the bedroom with a sword.” There are three prepositional phrases: “at midnight”, “in the bedroom,” and “with a sword.” Each of these can either be describing the main action (e.g. the killing was done with a sword) or the noun right before it (e.g. there are several bedrooms and we mean the one with a sword hanging over the dresser). Since each decision is independent, a sentence with k prepositional phrases like this has 2^k possible parses.¹

So it's really bad to organize parsing by considering all possible parse trees. Instead, consider all substrings in our input. If w has length n , then it has $\sum_{k=1}^n (n - k + 1) = \frac{n(n+1)}{2} = \binom{n+1}{2}$ substrings.² CYK computes a table summarizing the possible parses for each substring. From the table, we can quickly tell whether an input has a parse and extract one representative parse tree.³

1.2 CYK by example

Suppose the input sentence w is “Jeff trains geometry students” and the grammar has start symbol S and the following rules:

¹In real life, long sentences in news reports often exhibit versions of this problem.

²Draw an n by $n + 1$ rectangle and fill in the lower half.

³It still takes exponential time to extract all parse trees from the table, but we usually interested only in one of these trees.

\Rightarrow	$S \rightarrow N V_P$
	$VN \rightarrow N N$
	$V_P \rightarrow V N$
	$N \rightarrow \text{students} \mid \text{Jeff} \mid \text{geometry} \mid \text{trains}$
	$V \rightarrow \text{trains}$

Given a string w of length n , we build a triangular table with n rows and n columns. Conceptually, we write w below the bottom row of the table. The i th column correspond to the i th word. The cell at the i th column and the j th row (from the bottom) of the table corresponds to the substring starting the i th character of length j . The following is the table, and the substrings each entry corresponds to.

len				
4	Jeff trains geometry students			
3	Jeff trains geometry	trains geometry students		
2	Jeff trains	trains geometry	geometry students	
1	Jeff	trains	geometry	students
	Jeff	trains	geometry	students

first word in substring

CYK builds a table containing a cell for each substring. The cell for a substring x contains a list of variables V from which we can derive x (in one or more steps).

length	4				
	3				
	2				
	1				
		Jeff	trains	geometry	students

first word in substring

The bottom row contains the variables that can derive each substring of length 1. This is easy to fill in:

length	4				
	3				
	2				
	1	N	N,V	N	N
		Jeff	trains	geometry	students

first word in substring

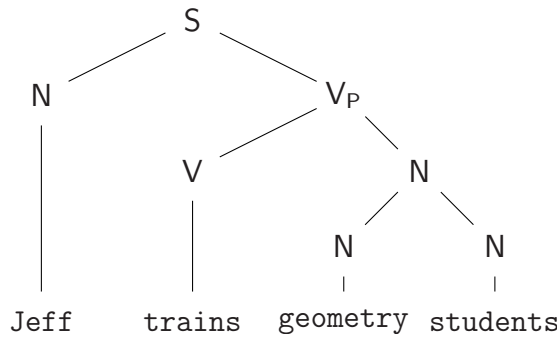
Now we fill the table row-by-row, moving upwards. To fill in the cell for a 2-word substring x , we look at the labels in the cells for its two constituent words and see what rules could derive this pair of labels. In this case, we use the rules $N \rightarrow N N$ and $V_P \rightarrow V N$ to produce:


```

CYK (  $\mathcal{G}, w$  )
     $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S), \Sigma \cup \mathcal{V} = \{X_1, \dots, X_r\}, w = w_1 w_2 \dots w_n.$ 
begin
    Initialize the 3d array  $B[1 \dots n, 1 \dots n, 1 \dots r]$  to FALSE
    for  $i = 1$  to  $n$  do
        for  $(X_j \rightarrow x) \in \mathcal{R}$  do
            if  $x = w_i$  then  $B[i, i, j] \leftarrow \text{TRUE}$ .
    for  $i = 2$  to  $n$  do /* Length of span */
        for  $L = 1$  to  $n - i + 1$  do /* Start of span */
             $R = L + i - 1$  /* Current span  $s = w_L w_{L+1} \dots w_R$  */
            for  $M = L + 1$  to  $R$  do /* Partition of span */
                /*  $x = w_L w_{L+1} \dots w_{M-1}, y = w_M w_{M+1} \dots w_R$ , and  $s = xy$  */
                for  $(X_\alpha \rightarrow X_\beta X_\gamma) \in \mathcal{R}$  do
                    /* Can we match  $X_\beta$  to  $x$  and  $X_\gamma$  to  $y$ ? */
                    if  $B[L, M - 1, \beta]$  and  $B[M, R, \gamma]$  then
                         $B[L, R, \alpha] \leftarrow \text{TRUE}$  /* If so, then can generate  $s$  by  $X_\alpha$ ! */
    for  $i = 1$  to  $r$  do
        if  $B[1, n, i]$  then return TRUE
    return FALSE

```

Figure 1: The CYK algorithm.



We have $O(n^2)$ cells in the table. For each cell, we have to consider n ways to divide its substring into two smaller substrings. So the table-filling procedure takes only $O(n^3)$ time.

2 The CYK algorithm

In general, we get the following result.

Theorem 2.1 *Let $\mathcal{G} = (\mathcal{V}, \Sigma, \mathcal{R}, S)$ be a grammar in CNF with $r = |\Sigma| + |\mathcal{V}|$ variables and terminals, and $t = |\mathcal{R}|$ rules. Let $w \in \Sigma^*$ be a word of length n . Then, one can compute a parse tree for w using \mathcal{G} , if $w \in L(\mathcal{G})$. The running time of the algorithm is $O(n^3 t)$.*

The result just follow from the CYK algorithm depicted in Figure 1. Note, that our pseudo-code just decides if a word can be generated by a grammar. With slight modifications, one can even generate the parse tree.