

CS344 - CYK Algorithm + Parsing

Note Title

2/3/2012

Announcements

- HW2 due now

- HW3 is up

↳ feel free to use the web

due on Sunday

WHENEVER I LEARN A NEW SKILL I CONCOCT ELABORATE FANTASY SCENARIOS WHERE IT LETS ME SAVE THE DAY.

OH NO! THE KILLER MUST HAVE FOLLOWED HER ON VACATION!



BUT TO FIND THEM WE'D HAVE TO SEARCH THROUGH 200 MB OF EMAILS LOOKING FOR SOMETHING FORMATTED LIKE AN ADDRESS!

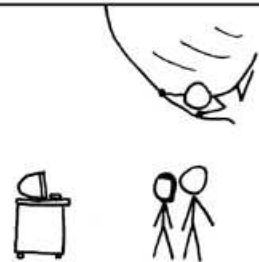


IT'S HOPELESS!

EVERYBODY STAND BACK.



I KNOW REGULAR EXPRESSIONS.



Recap: CFGs and CNF

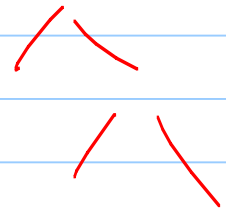
Chomsky Normal Form:

Each rule in the grammar is either:

- $A \rightarrow BC$ where BC are 2 nonterminals
where neither B or C is the start variable, & both are nonterminals
- $A \rightarrow a$ where a is a terminal
- No useless symbols

Why do we care?

① Makes structure of parsing nice.



② Computable: $O(n^2)$ time to convert to CNF

③ Given CNF, can compute a parse tree in $O(n^3)$ time (using dynamic programming).

Ex: Simple CFG

$S_0 \rightarrow \underline{A}B$

$A \rightarrow aA \mid \epsilon$

$B \rightarrow bA \underline{bAB} \mid \epsilon$

$a^* (b a^* b a^*)^*$

tokens or terminals

Convert:

① useless states?

No

② Eliminate ϵ -transitions

$$S_0 \rightarrow AB$$

$$A \rightarrow aA \quad | \quad \cancel{\epsilon}$$

$$B \rightarrow bAbAB \quad | \quad \cancel{\epsilon}$$

↪

$$S_0 \rightarrow AB \quad | \quad A \quad | \quad B \quad | \quad \epsilon$$

$$A \rightarrow aA \quad | \quad a$$

$$B \rightarrow bAbAB \quad | \quad bAbA \quad | \quad bb \quad | \quad bbB$$
$$bAb \quad | \quad bbA \quad | \quad bbAB \quad | \quad bAbB$$

$X \rightarrow Y$
 $X \rightarrow Y \rightarrow Z$

③ Eliminate unit pairs : $(S, A) + (S, B)$

$S_0 \rightarrow AB \mid A \mid B \mid \epsilon$

$A \rightarrow aA \mid a$

$B \rightarrow bAbAB \mid bAbA \mid bb \mid bbB \mid bAb \mid bbA \mid bbAB \mid bAbB$

$S_0 \rightarrow AB \mid aA \mid a \mid bAbAB \mid bAbA \mid bb \mid bbB \mid bAb \mid bbA \mid bbAB \mid bAbB \mid \epsilon$
 $A \rightarrow aA \mid a$
 $B \rightarrow bAbAB \mid bAbA \mid bb \mid bbB \mid bAb \mid bbA \mid bbAB \mid bAbB$

④ a: Get rid of non-single terminals
 (eg $X \rightarrow yz$ becomes $X \rightarrow YZ$ & $Y \rightarrow y$)

$S_0 \rightarrow$ $\begin{matrix} AB & | & aA & | & a \\ bAb & | & bAbA & | & bb & | & bbB \\ bAb & | & bbA & | & bbAB & | & bAbB \end{matrix} \quad | \quad \epsilon$
 $A \rightarrow$ $\begin{matrix} XA & | & a \\ YAY & | & AB \\ bAb & | & bbA & | & bbAB & | & bAbB \end{matrix}$
 $B \rightarrow$ $\begin{matrix} YAY & | & AB \\ bAb & | & bbA & | & bbAB & | & bAbB \end{matrix}$

$\begin{matrix} X \rightarrow a \\ Y \rightarrow b \end{matrix}$

(replace any a with X
 & any b with Y)

4b): Eliminate 3 or more non-terminal transitions

$B \rightarrow Y A Y A B$ (too long)

$\left\{ \begin{array}{l} B \rightarrow Y M_1 \\ M_1 \rightarrow A M_2 \\ M_2 \rightarrow Y M_3 \\ M_3 \rightarrow A B \end{array} \right.$

$B \rightarrow M_1 M_2$

$M_1 \rightarrow Y A$

$M_2 \rightarrow M_1 B$

Now, parsing. Consider
Is language?

abbaabab.

How to parse:

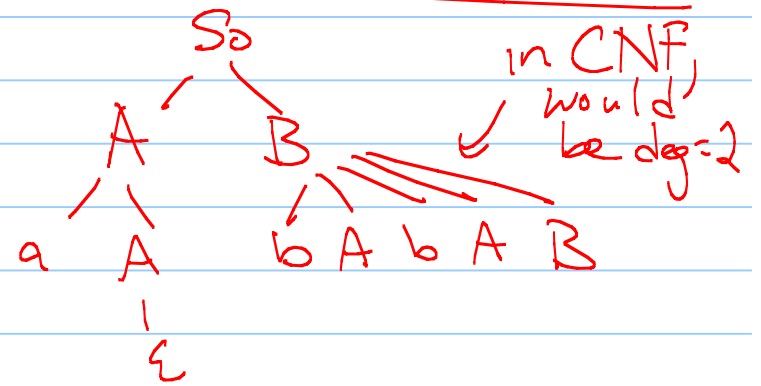
$S_0 \rightarrow AB$

$A \rightarrow aA \mid \epsilon$

$B \rightarrow bAbAB \mid \epsilon$

tokens

① $S_0 \xrightarrow{(do A)} AB \rightarrow aAB$
 $\rightarrow aB$
 $\rightarrow abAB$
 $\rightarrow abbAB$



Parse tree: (prev page)

CKK algorithm: build a table

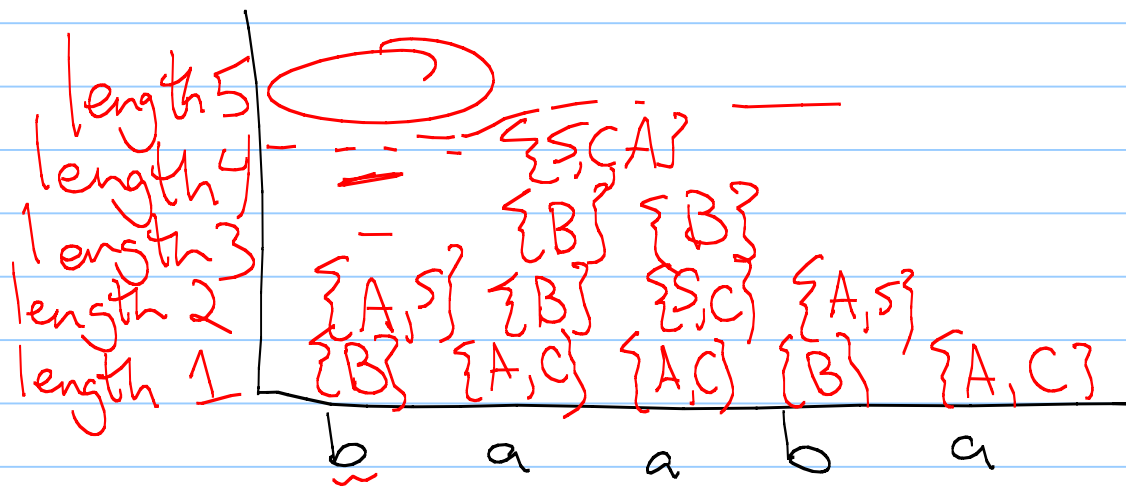
Given a word $w = w_1 w_2 w_3 w_4 \dots w_k$,
we'll look at all possible
substrings $w_i w_{i+1} \dots w_{j-1} w_j$,
and look at how they can be
parsed.

We'll build a table from the bottom up.

Ex: $S \rightarrow AB \mid BC$
 $A \rightarrow BA \mid a$
 $B \rightarrow CC \mid b$
 $C \rightarrow AB \mid a$

$X \rightarrow x$
 Complexity \hookrightarrow

Compute valid parse tree for 'baaba'



Running times:

Say we have n rules.

Converting to CNF: $O(n^2)$

finding unit pairs

Running CYK: $O(n^3)$

Other parsing algorithms

CYK is still pretty slow, especially for large programming languages.

After it was developed, a lot of work was put into figuring out what grammars could have faster algorithms.

Two big (& useful) classes have linear time parsers: LL & LR.