

CS 3200 - Parsing

Note Title

1/25/2016

Announcements

- HW due Saturday

Last time - Ch 2, Sec. 3

- Context free languages
- Automated parsing tools
 - bison

Examples: Palindromes

Recursive defn: First & last letters match & interior is a palindrome.
alphabet {a, b}

$S \rightarrow aSa \mid bSb \mid \epsilon$

← or

$S \rightarrow a \mid b$

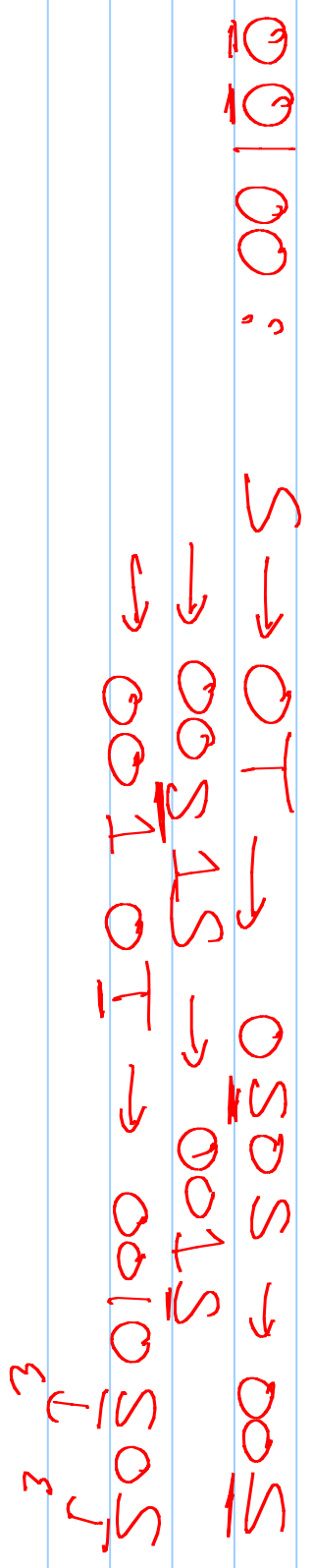
Ex: Derive a CFG for all binary strings w/ an even # of 0's.

S \rightarrow S0S0S | ϵ | S1S -

\uparrow

0011

An alternate for that last one:



Context-Free Languages

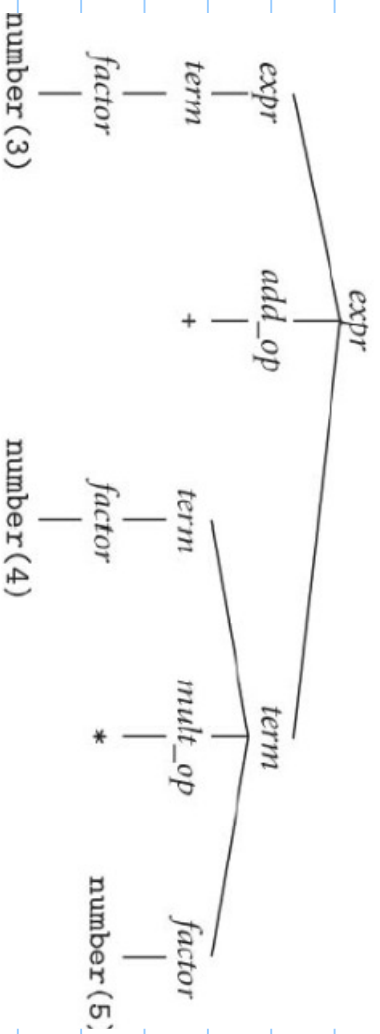
Recall that for any context free language there are an infinite # of grammars that can produce it.

We wish to somehow give a definition of a 'good' set of productions.

Goal: Parsing (well) -
given a language, detect if
a string is in that
language.

Why?

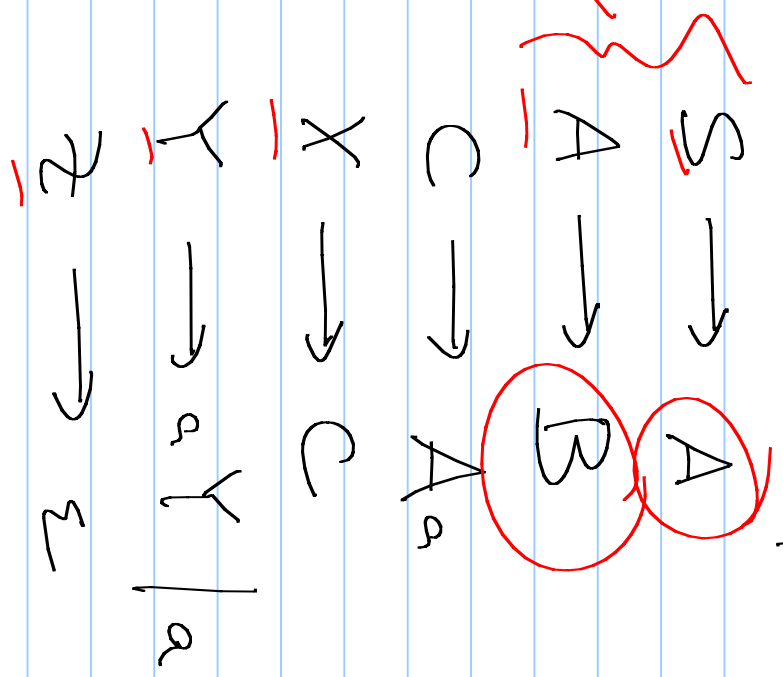
- Remember, goal is to compile a program. \downarrow
- Need to be able to quickly check if it is valid and give a parse tree back!



So - a bad example: $S_0 \rightarrow S | X | Z$

Why is it bad?

- useless states
- unreachable non-terminals
- wait pair



Chomsky Normal Form (CNF)

Each rule in the grammar is either:

- $A \rightarrow BC$

where neither B or C is the start variable, or both are non-terminals

- $A \rightarrow a$

where a is a terminal

- $S \rightarrow \epsilon$

where S is the start symbol

But: Clearly, not always in CNF!

Luckily:

Thm: All grammars can be converted to CNF.

Procedure:

① Create a new start symbol S_0 ,
a $S_0 \rightarrow S$

② Eliminate useless rules

(just delete ones that
can't be reached)

② Remove nullable variables.
 $A \rightarrow \epsilon$

How?

Remove all ϵ productions.
Then fix.

Ex: $A \rightarrow CB C$ } $A \rightarrow CB C | CC$

~~$B \rightarrow \epsilon$~~
 $B \rightarrow B$

$D \rightarrow B A C B$ } $D \rightarrow B A C B | A C B |$
 $B A C | A C$

③ Remove unit rules:

$$S \rightarrow A$$

How? Must have:

$$X \rightarrow z_1, z_1 \rightarrow z_2, \dots, z_n \rightarrow Y$$

(Since we removed ϵ -transitions in ②)

Then:

$$\begin{array}{l} S \rightarrow XAY \\ \underbrace{A \rightarrow B|XY} \\ B \rightarrow b|c \end{array} \rightarrow \begin{array}{l} S \rightarrow XbY|XcY \\ A \rightarrow XY \end{array}$$

(Skill ③)

For each unit pair (A, B)
and rule $B \rightarrow w$,

add $A \rightarrow w$ to a new
grammar.

Note that (A, A) is a unit pair,
so all rules $A \rightarrow w$
will stick around.)

④ Get rid of "long" right hand sides.

4a: Create $V_c \rightarrow c$ for every character.

Replace c with V_c everywhere.

Now either

$A \rightarrow CDEF$

or

$V_c \rightarrow c$.

4b: $A \rightarrow B_1 B_2 B_3 \dots B_k$

How to replace with only 2 nonterminals on the right?

$A \rightarrow B_1 X_1$
 $X_1 \rightarrow B_2 X_2$
 $X_2 \rightarrow B_3 X_3$
 \vdots
 $X_{k-2} \rightarrow B_{k-1} B_k$

$A \rightarrow W X Y Z$
 $A \rightarrow W D B$
 $B \rightarrow X C$
 $C \rightarrow Y Z$

Ex:

$S \rightarrow ASA \mid aB$

$A \rightarrow B \mid S$

$B \rightarrow b \mid \epsilon$

Now - why do we care??

Parsing: building those parse trees
we saw

In general, there are an exponential
number of parse trees
for a given input.

So how to check quickly?

Even in CNF might be 2ⁿ possible
parse trees.

Cocke - Younger - Kasami (CYK) algorithm

Uses a table & dynamic programming
to give a parse tree in $O(n^3)$ time.

Grammar must be in CNF!

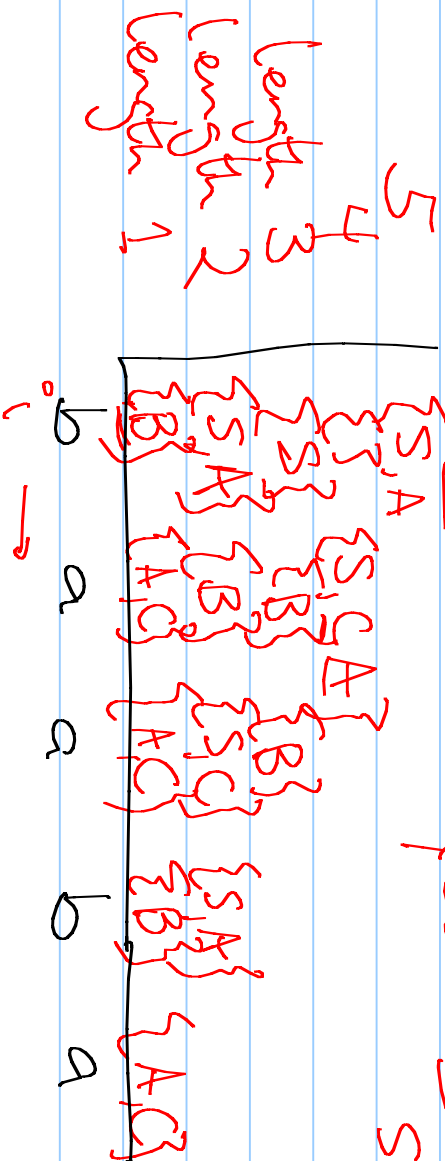
CYK Algorithm

Given a word $w = w_1 w_2 w_3 \dots w_k$
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.

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

Test if 'babbe' is in the language
 yes — S is in top spot



Running time:

Say we have n rules.

Converting to CNF:

$O(n^2)$

Running CYK: n^3 time, n^2 space

Other parsing algorithms

CRK 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 (or useful) classes have linear time parsers: LL & LR.

