

MACM 300

Formal Languages and Automata

Anoop Sarkar

<http://www.cs.sfu.ca/~anoop>

Derivations of a CFG

- *strings grow on trees*
- *strings grow on Noun*
- *strings grow Object*
- *strings Verb Object*
- **Noun Verb Object**
- **Sentence**

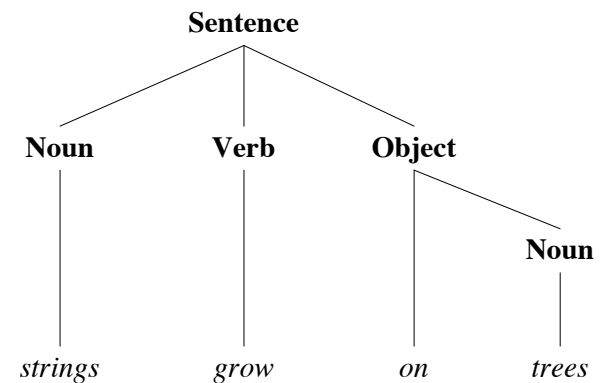
3

Context-free Grammars

- Set of rules by which valid sentences can be constructed.
- Example:
 - Sentence \rightarrow Noun Verb Object
 - Noun \rightarrow *trees* | *strings*
 - Verb \rightarrow *are* | *grow*
 - Object \rightarrow *on* Noun | Adjective
 - Adjective \rightarrow *slowly* | *interesting*
- What strings can Sentence *derive*?
- Syntax only – no semantic checking

2

Derivations and parse trees



4

CFG Notation

- Normal CFG notation

$$E \rightarrow E * E$$

$$E \rightarrow E + E$$

- Backus Naur notation

$$E ::= E * E \mid E + E$$

(an or-list of right hand sides)

5

CFG: Formal Definition

- A CFG is a 4-tuple (V, Σ, R, S) , where:
 - V is a finite non-empty set of symbols (called *variables*, or *non-terminals*)
 - Σ is a finite set of symbols (called set of terminal symbols, or the alphabet)
 - We generally assume that $V \cap \Sigma = \emptyset$
 - R is a finite non-empty set of rules, where each rule is of the form: $A \rightarrow w$, $w \in (V \cup \Sigma)^*$
 - $S \in V$ is the start variable (or start non-terminal)

6

CFGs and Languages

- For the regular expression $(01)^*$ we know that it corresponds to a language $L = \{\epsilon, 01, 0101, 010101, \dots\}$
- We know this because we know the definition of the Kleene closure operator $*$
- Similarly we define a relation between CFGs and a language as follows:
 1. Write down the start variable
 2. Find a variable written down so far and a rule that has that variable on the left-hand side of a rule in the CFG. Replace the variable with the right-hand side of the rule
 3. Repeat step 2 until no variable remains

7

CFGs and Languages

- We can formally specify the 3 rules for deriving a language from a CFG as follows:
 - If u, v, w are strings from $(V \cup \Sigma)^*$
 - Then if we have a string uAv written down
 - And if we have a rule in the CFG, $A \rightarrow w$,
 - Then we can replace A with w , written as:
 $uAv \Rightarrow wwv$ (called a **derivation** step)

8

CFGs and Languages

- Consider CFG: $A \rightarrow 0A1 \mid \epsilon$, here are some strings derived from this grammar:
 - $A \Rightarrow \epsilon$
 - $A \Rightarrow 0A1 \Rightarrow 0\epsilon 1 = 01$ (derived in 0 or more steps: $A \Rightarrow^* 01$)
 - $A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 00\epsilon 11 = 0011$
 - $A \Rightarrow 0A1 \Rightarrow 00A11 \Rightarrow 000A111 \Rightarrow 000\epsilon 111 = 000111$
 - ...and so on...
- The set of all strings that can be derived is the *language* of the CFG
 - In this case the language is $\{ 0^n 1^n \mid n \geq 0 \}$

9

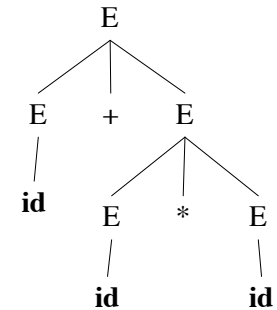
Arithmetic Expressions

- $E \rightarrow E + E$
- $E \rightarrow E * E$
- $E \rightarrow (E)$
- $E \rightarrow - E$
- $E \rightarrow \text{id}$

10

Leftmost derivations for $\text{id} + \text{id} * \text{id}$

- $E \Rightarrow E + E$
- $\Rightarrow \text{id} + E$
- $\Rightarrow \text{id} + E * E$
- $\Rightarrow \text{id} + \text{id} * E$
- $\Rightarrow \text{id} + \text{id} * \text{id}$

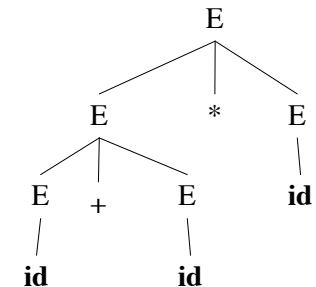


$E \Rightarrow_{\text{lm}}^* \text{id} + \text{id} * E$

11

Leftmost derivations for $\text{id} + \text{id} * \text{id}$

- $E \Rightarrow E * E$
- $\Rightarrow E + E * E$
- $\Rightarrow \text{id} + E * E$
- $\Rightarrow \text{id} + \text{id} * E$
- $\Rightarrow \text{id} + \text{id} * \text{id}$



12

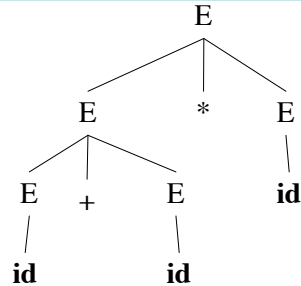
Rightmost derivation for

id + id

Note that the parse tree is the same as one of the leftmost derivations.

$E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow (E)$
 $E \rightarrow - E$
 $E \rightarrow id$

$E \Rightarrow E * E$
 $\Rightarrow E * id$
 $\Rightarrow E + E * id$
 $\Rightarrow E + id * id$
 $\Rightarrow id + id * id$



$E \Rightarrow_{rm}^* E + E * id$

A grammar G is **ambiguous** iff there are **two parse trees** or two distinct leftmost or two distinct rightmost derivations for **any** string in L(G)

13

Ambiguity

- Alternatives
 - Massage grammar to make it unambiguous
 - Rely on “default” parser behavior
 - Augment parser
- Consider the original ambiguous grammar:
 - $E \rightarrow E + E$ $E \rightarrow E * E$
 - $E \rightarrow (E)$ $E \rightarrow - E$
 - $E \rightarrow id$
- How can we change the grammar to get only one tree for the input **id + id * id**

15

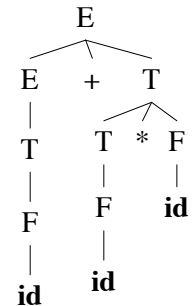
Ambiguity

- Grammar is ambiguous if more than one parse tree is possible for some sentences
- Examples in English:
 - Two sisters reunited after 18 years in checkout counter
- Ambiguity is not acceptable in PL
 - Unfortunately, it’s *undecidable* to check whether a grammar is ambiguous

14

Ambiguity

- Original ambiguous grammar:
 - $E \rightarrow E + E$ $E \rightarrow E * E$
 - $E \rightarrow (E)$ $E \rightarrow - E$
 - $E \rightarrow id$
- Unambiguous grammar:
 - $E \rightarrow E + T$ $T \rightarrow T * F$
 - $E \rightarrow T$ $T \rightarrow F$
 - $F \rightarrow (E)$ $F \rightarrow - E$
 - $F \rightarrow id$
- Input: id + id * id



16

Other Ambiguous Grammars

- Consider the grammar
 $R \rightarrow R \cup R \mid RR \mid R^* \mid ('R')$
- What does this grammar generate?
- What's the parse tree for $a \cup b^* a$
- Is this grammar ambiguous?

17

Chomsky Normal Form

- First step, add S_0 and then remove epsilon rules
 $S \rightarrow ASA \mid aB$
 $A \rightarrow B \mid S$
 $B \rightarrow b \mid \epsilon$
- After adding S_0 and then ϵ -removal (remove $*$):
 $S_0 \rightarrow S$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$
 $A \rightarrow B \mid S \mid \epsilon^*$
 $B \rightarrow b \mid \epsilon^*$

19

Grammar Transformations

- G is converted to G' s.t. $L(G') = L(G)$
- Left Factoring
- Removing cycles: $A \Rightarrow^+ A$
- Removing ϵ -rules of the form $A \rightarrow \epsilon$
- Eliminating left recursion
- Conversion to normal forms:
 - Chomsky Normal Form, $A \rightarrow BC$ and $A \rightarrow a$
 - Greibach Normal Form, $A \rightarrow a\beta$

18

Chomsky Normal Form

- Second step, remove unit rules or chain rules
 $S_0 \rightarrow S$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$
 $A \rightarrow B \mid S$
 $B \rightarrow b$
- After removal of unit rules $S_0 \rightarrow S$ and $S \rightarrow S$:
 $S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $A \rightarrow B \mid S$
 $B \rightarrow b$

20

Chomsky Normal Form

- After removal of unit rules $S_0 \rightarrow S$ and $S \rightarrow S$:
 $S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $A \rightarrow B \mid S$
 $B \rightarrow b$
- After removal of unit rules $A \rightarrow B$ and $A \rightarrow S$:
 $S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $A \rightarrow \mathbf{b} \mid \mathbf{ASA} \mid \mathbf{aB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$
 $B \rightarrow b$

21

Chomsky Normal Form

- Converting the rhs of each rule to have two non-terminals
 $A \rightarrow B N_1 C N_2$
 $N_1 \rightarrow a$
 $N_2 \rightarrow d$
- After converting to binary form:
 $A \rightarrow B N_3 \quad N_1 \rightarrow a$
 $N_3 \rightarrow N_1 N_4 \quad N_2 \rightarrow d$
 $N_4 \rightarrow C N_2$

23

Chomsky Normal Form

Consider some simpler examples for next two steps

- Removing terminals from the rhs of rules
 $A \rightarrow B a C d$
- After removal of terminals from the rhs:
 $A \rightarrow B N_1 C N_2$
 $N_1 \rightarrow a$
 $N_2 \rightarrow d$

22

Chomsky Normal Form

Back to original example

- After removal of unit rules $A \rightarrow B$ and $A \rightarrow S$:
 $S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$
 $A \rightarrow \mathbf{b} \mid \mathbf{ASA} \mid \mathbf{aB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$
 $B \rightarrow b$
- After adding variables:
 $S_0 \rightarrow \mathbf{AA}_1 \mid \mathbf{UB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$
 $S \rightarrow \mathbf{AA}_1 \mid \mathbf{UB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$
 $A \rightarrow \mathbf{b} \mid \mathbf{AA}_1 \mid \mathbf{UB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$
 $\mathbf{A}_1 \rightarrow \mathbf{SA} \quad \mathbf{U} \rightarrow \mathbf{a} \quad \mathbf{B} \rightarrow \mathbf{b}$

24

Non-CF Languages

- The pumping lemma for CFLs [Bar-Hillel] is similar to the pumping lemma for RLs
- For a string $wuxvy$ in a CFL for $u, v \neq \varepsilon$ and the string is long enough then $wu^n xv^n y$ is also in the CFL for $n \geq 0$
- Not strong enough to work for every non-CF language (cf. Ogden's Lemma)

25

Non-CF Languages

$$L_1 = \{wcbw \mid w \in (a|b)^*\}$$

$$L_2 = \{a^n b^m c^n d^m \mid n \geq 1, m \geq 1\}$$

$$L_3 = \{a^n b^n c^n \mid n \geq 0\}$$

26

CF Languages

$$L_4 = \{wcbw^R \mid w \in (a|b)^*\}$$

$$S \rightarrow aSa \mid bSb \mid c$$

$$L_5 = \{a^n b^m c^m d^n \mid n \geq 1, m \geq 1\}$$

$$S \rightarrow aSd \mid aAd$$

$$A \rightarrow bAc \mid bc$$

27

Context-free languages and Pushdown Automata

- Recall that for each regular language there was an equivalent finite-state automaton
- The FSA was used as a recognizer of the regular language
- For each context-free language there is also an automaton that recognizes it: called a **pushdown automaton (pda)**

28

Context-free languages and Pushdown Automata

- Similar to FSAs there are non-deterministic pda and deterministic pda
- Unlike in the case of FSAs we cannot always convert a npda to a dpda
- Similar to the FSA case, a DFA construction provides us with the algorithm for lexical analysis,
- In this case the construction of a dpda will provide us with the algorithm for parsing (take in strings and provide the parse tree)

29

Summary

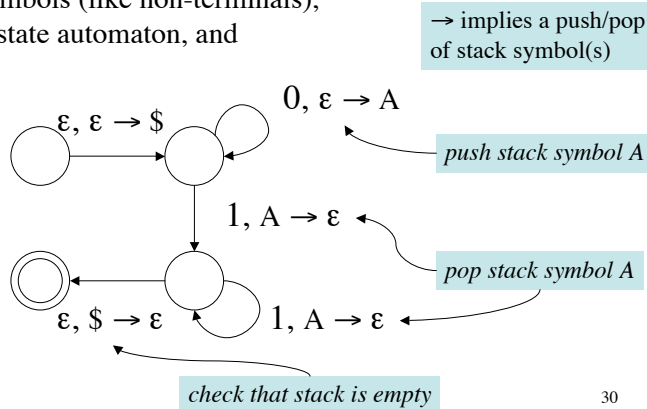
- CFGs can be used describe PL
- Derivations correspond to parse trees
- Parse trees represent structure of programs
- Ambiguous CFGs exist
- CF languages can be recognized using Pushdown Automata

31

Pushdown Automata

- PDA has
 - an alphabet (terminals) and
 - stack symbols (like non-terminals),
 - a finite-state automaton, and
 - stack

e.g. PDA for language
 $L = \{ 0^n 1^n : n \geq 0 \}$



30