

Context-Free Grammars

Sipser 2.1 (pages 99 – 109)

What are we missing?

- So far:
- We know how to *recognize* languages
 - With finite state automata
 - As people...
- We know how to *generate* languages
 - With regular expressions
 - As people...
- Finite state automata and regular expressions are limited, though!

Bring back memories?

- In English, a grammar tells us whether a particular sentence is well formed or not
- For instance, “a sentence can consist of a noun phrase followed by a verb phrase”
- More concisely, we could write
 $\langle sentence \rangle \rightarrow_G \langle noun_phrase \rangle \langle verb_phrase \rangle$

Great, but what's a noun phrase?

- A sentence is
 - $\langle \text{sentence} \rangle \rightarrow_G \langle \text{noun_phrase} \rangle \langle \text{verb_phrase} \rangle$
- We need to provide definitions for the newly introduced constructs
 - $\langle \text{noun_phrase} \rangle$ and $\langle \text{verb_phrase} \rangle$
 - $\langle \text{noun_phrase} \rangle \rightarrow_G \langle \text{article} \rangle \langle \text{noun} \rangle$
 - $\langle \text{verb_phrase} \rangle \rightarrow_G \langle \text{verb} \rangle$

Generating well-formed sentences

- Grammar rules so far:
 - $\langle \text{sentence} \rangle \rightarrow_G \langle \text{noun_phrase} \rangle \langle \text{verb_phrase} \rangle$
 - $\langle \text{noun_phrase} \rangle \rightarrow_G \langle \text{article} \rangle \langle \text{noun} \rangle$
 - $\langle \text{verb_phrase} \rangle \rightarrow_G \langle \text{verb} \rangle$
- *To complete our simple grammar, we associate actual words with the terms $\langle \text{article} \rangle$, $\langle \text{noun} \rangle$, and $\langle \text{verb} \rangle$*
 - $\langle \text{article} \rangle \rightarrow_G a$
 - $\langle \text{article} \rangle \rightarrow_G the$
 - $\langle \text{noun} \rangle \rightarrow_G student$
 - $\langle \text{verb} \rangle \rightarrow_G relaxes$
 - $\langle \text{verb} \rangle \rightarrow_G studies$

Context-free grammars

- A context-free grammar G is a quadruple (V, Σ, R, S) , where
 - V is a finite set called the **variables**
 - Σ is a finite set, disjoint from V , called the **terminals**
 - R is a finite subset of $V \times (V \cup \Sigma)^*$ called the **rules**
 - $S \in V$ is called the **start symbol**
- **For any $A \in V$ and $u \in (V \cup \Sigma)^*$,
we write $A \rightarrow_G u$ whenever $(A, u) \in R$**

The language of a grammar

- If
 - $u, v, w \in (V \cup \Sigma)^*$
 - $A \rightarrow_G w$ is a rulethen
 - We say uAv **yields** uwv
 - **Write** $uAv \Rightarrow_G uwv$
- **If**
 - $u \Rightarrow_G u1 \Rightarrow_G u2 \Rightarrow_G \dots \Rightarrow_G uk \Rightarrow_G v$**then**
 - **We write** $u \Rightarrow_G^* v$
- **The language of the grammar G is**
$$L(G) = \{w \in \Sigma^* \mid S \Rightarrow_G^* w\}$$

For example...

- Consider $G = (V, \Sigma, R, S)$, where
 - $V = \{S\}$
 - $\Sigma = \{a, b\}$
 - $R = \{ S \rightarrow_G aSa \mid bSb \mid aSb \mid bSa \mid \varepsilon \}$
- Is there a grammar whose language is $PAL = \{w \in \Sigma^* \mid w = reverse(w)\}$?

Arithmetic expressions and parse trees

- Consider $G = (V, \Sigma, R, S)$, where
 - $V = \{ \langle \text{EXPR} \rangle, \langle \text{TERM} \rangle, \langle \text{FACTOR} \rangle \}$
 - $\Sigma = \{ a, +, \times, (,) \}$
 - $R = \{ \langle \text{EXPR} \rangle \rightarrow_G \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle, \langle \text{TERM} \rangle \rightarrow_G \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle, \langle \text{FACTOR} \rangle \rightarrow_G \langle \text{EXPR} \rangle \mid a \}$
 - $S = \langle \text{EXPR} \rangle$
- What about $a \times a + a$?

Leftmost derivation

- A **derivation** of a string in a grammar is a **leftmost derivation** if:
 - at every step the *leftmost* remaining variable is the one replaced

Needlessly complicated?

- How about just

$$\begin{aligned} \langle \text{EXPR} \rangle \rightarrow_G \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle \mid \\ \langle \text{EXPR} \rangle * \langle \text{EXPR} \rangle \mid \\ a \end{aligned}$$

- A grammar G is **ambiguous** if some string w has two or more different leftmost derivations

Chomsky normal form

- A context-free grammar G is in **Chomsky normal form**
 - If every rule is of the form
 - $A \rightarrow BC$
 - $A \rightarrow a$
 - **where** $A, B, C \in V$, $B \neq S \neq C$, **and** $a \in \Sigma$
 - **We permit $S \rightarrow \varepsilon$**

Chomsky normal form

- **Theorem 2.9: Any context-free language is generated by a context-free grammar in Chomsky normal form**
- **Proof:**
 1. **Make sure S appears only on the left**
 2. **Remove empty rules: $A \rightarrow \varepsilon$**
 3. **Handle unit rules: $A \rightarrow B$**
 4. **Fix all the rest...**
- **For example:**
 - $S \rightarrow_G ASA \mid aA$
 - $A \rightarrow_G b \mid \varepsilon$