

5 Context-Free Grammars

5.1 Context-free grammars

Reading: Sipser 2.1 (pp. 100-105)

- CFG $G = (V, \Sigma, R, S)$, where V is the set of variables, Σ is the set of terminals (alphabet), R is the set of rules in the form of $V \rightarrow (V \cup \Sigma)^*$ (head \rightarrow body), and $S \in V$ is the start variable.
- The CFG that generates all palindromes (strings that read the same forward and backward) over $\{0, 1\}$ is $G = (\{S\}, \{0, 1\}, R, S)$, where R contains $S \rightarrow \epsilon | 0|1|0S0|1S1$.
- Let u, v, w be strings in $(V \cup \Sigma)^*$. If $A \rightarrow w$ is a rule, then uAv yields uwv , written $uAv \Rightarrow uwv$. We say u derives v , written $u \xRightarrow{*} v$, if $\exists u_1, \dots, u_k \in (V \cup \Sigma)^*$ such that $u \Rightarrow u_1 \Rightarrow \dots \Rightarrow u_k \Rightarrow v$. Here, \Rightarrow means one step and $\xRightarrow{*}$ means zero or more steps.
- Leftmost and rightmost derivations: $\xRightarrow{*}_{lm}$, $\xRightarrow{*}_{lm}$, $\xRightarrow{*}_{rm}$, $\xRightarrow{*}_{rm}$.
- The language of a CFG G , $L(G) = \{w \in \Sigma^* | S \xRightarrow{*} w\}$. $L(G)$ is said to be a CFL.
- Why called “context-free”?

Example: Is $L = \{0^n 1^n | n \geq 0\}$ a context-free language? Yes since it is the language of the context-free grammar $S \rightarrow 0S1 | \epsilon$.

Example: What is the language for grammar G with $S \rightarrow AA$ and $A \rightarrow AAA | bA | Ab | a$? $L(G) = \{w \in \{a, b\}^* | w \text{ has an even (nonzero) number of } a\}$.

Example: A CFG for simple expressions in programming languages: $E \rightarrow E + E | E * E | (E) | I$ and $I \rightarrow Ia | Ib | I0 | I1 | a | b$.

5.2 Parse trees

- A parse tree is a tree representation for a derivation, in which each interior node is a variable, each leaf node is either a variable, a terminal, or ϵ , and if an interior node is a variable A and its children are X_1, \dots, X_k , then there must be a rule $A \rightarrow X_1 \dots X_k$.
- Yield of a parse tree: Concatenation of the leaf nodes in a parse tree rooted at the start variable.
- Four equivalent notions:
 1. $A \xRightarrow{*} w$;
 2. $A \xRightarrow{*}_{lm} w$;
 3. $A \xRightarrow{*}_{rm} w$; and
 4. A parse tree with root A and yield w .

5.3 Regular grammars

- A regular grammar (RG) is a special CFG where the body of each rule contains at most one variable which, if present, must be the last symbol in the body.
- For each regular language L there is a RG G such that $L = L(G)$.
- Given RG $G = (V, \Sigma, R, S)$, construct an equivalent NFA $N = (Q, \Sigma, \delta, q_0, F)$, where $Q = V \cup \{f\}$, $q_0 = S$, $F = \{f\}$, and for each rule $A \rightarrow wB$, $\delta(A, w) = B$ and for each rule $A \rightarrow w$, $\delta(A, w) = f$.
- Given DFA $A = (Q, \Sigma, \delta, q_0, F)$, define an equivalent RG $G = (V, \Sigma, R, S)$, where $V = Q$, $S = q_0$, and for each transition $\delta(q, a) = p$, $q \rightarrow ap$ and for each final state $q \in F$, $q \rightarrow \epsilon$.

Example: Given RG: $S \rightarrow aS | bA | aB$, $A \rightarrow aS | \epsilon$, and $B \rightarrow bS | \epsilon$, what is its NFA?

Example: Given a four-state DFA for language aa^*bb^*a , what is its regular grammar?

5.4 Ambiguity in grammars and languages

Reading: Sipser 2.1 (pp. 105-106)

- A CFG $G = (V, \Sigma, R, S)$ is ambiguous if there is $w \in \Sigma^*$ for which there are at least two parse trees (or leftmost derivations).
- Grammar $G: E \rightarrow E + E | E * E | (E) | I$ and $I \rightarrow Ia | Ib | IO | I1 | a | b$ is ambiguous since $a+b*a$ has two parse trees.
- Some ambiguous grammars have an equivalent unambiguous grammar. For example, an unambiguous grammar for the simple expressions is $G': E \rightarrow E + T | T, T \rightarrow T * F | F, F \rightarrow (E) | I$, and $I \rightarrow Ia | Ib | IO | I1 | a | b$.
- A context-free language is said to be inherently ambiguous if all its grammars are ambiguous. For example, $\{a^n b^n c^m d^m | m, n \geq 1\} \cup \{a^n b^m c^m d^n | m, n \geq 1\}$. Its grammar is $S \rightarrow S_1 | S_2, S_1 \rightarrow AB, A \rightarrow aAb | ab, B \rightarrow cBd | cd, S_2 \rightarrow aS_2d | aCd$, and $C \rightarrow bCc | bc$. The grammar is ambiguous (considering $abcd$). There is a not so easy proof that all grammars for the language are ambiguous, thus it is inherently ambiguous.
- There is no algorithm to determine whether a given CFG is ambiguous. There is no algorithm to remove ambiguity from an ambiguous CFG. There is no algorithm to determine whether a given CFL is inherently ambiguous.

5.5 Chomsky normal form

Reading: Sipser 2.1 (pp. 106-109)

For a CFL, different people may come up with different but equivalent CFGs. There are several steps to simplify a CFG.

- Eliminating ϵ rules for $L(G) - \{\epsilon\}$ by finding nullable variables ($A \xRightarrow{*} \epsilon$). For example, $S \rightarrow AB, A \rightarrow aAA | \epsilon$, and $B \rightarrow bBB | \epsilon$ can be changed to $S \rightarrow AB | A | B, A \rightarrow aAA | aA | a$, and $B \rightarrow bBB | bB | b$.
- Eliminating unit rules. For example, $E \rightarrow T | E + T, T \rightarrow F | T * F, F \rightarrow I | (E)$, and $I \rightarrow Ia | Ib | IO | I1 | a | b$ can be changed to $E \rightarrow E + T | T * F | (E) | Ia | Ib | IO | I1 | a | b, T \rightarrow T * F | (E) | Ia | Ib | IO | I1 | a | b, F \rightarrow (E) | Ia | Ib | IO | I1 | a | b$, and $I \rightarrow Ia | Ib | IO | I1 | a | b$.
- Eliminating useless variables (and thus associated rules) by finding nongenerating and unreachable variables. For example, $S \rightarrow AB | a$ and $A \rightarrow b$ can be simplified to $S \rightarrow a$.

The Chomsky Normal Form (CNF): Any nonempty CFL without ϵ has a CFG G in which all rules are in one of the following two forms: $A \rightarrow BC$ and $A \rightarrow a$, where A, B, C are variables, and a is a terminal. Note that one of the uses of CNF is to turn parse trees into binary trees.

To convert a CFG to a grammar in CNF:

- Add a new start variable S_0 in the case when the old start variable S appears in the body of some rules.
- Simplify the grammar by removing ϵ rules, unit rules, and useless variables.
- Convert the rules in the simplified grammar into the proper forms of CNF by adding additional variables and rules.

Example (Sipser p. 108): Given a CFG G , construct a CNF G' such that $L(G) - \{\epsilon\} = L(G')$.