

Context-Free Grammars and Languages

- We have seen that many languages cannot be regular. Thus we need to consider larger classes of langs.
- *Context-Free Languages* (CFL's) played a central role natural languages since the 1950's, and in compilers since the 1960's.
- *Context-Free Grammars* (CFG's) are the basis of BNF-syntax.
- Today CFL's are increasingly important for XML and their DTD's.

We'll look at: CFG's, the languages they generate, parse trees, pushdown automata, and closure properties of CFL's.

Informal example of CFG's

Consider $L_{pal} = \{w \in \Sigma^* : w = w^R\}$

For example otto $\in L_{pal}$, madamimadam $\in L_{pal}$.

In Finnish language e.g. saippuakauppias $\in L_{pal}$
(“soap-merchant”)

DoGeeseSeeGod
NoMelonNoLemon

Let $\Sigma = \{0, 1\}$ and suppose L_{pal} were regular.

Let n be given by the pumping lemma. Then $0^n 1 0^n \in L_{pal}$. In reading 0^n the FA must make a loop. Omit the loop; contradiction.

Let's define L_{pal} inductively:

Basis: ϵ , 0, and 1 are palindromes.

Induction: If w is a palindrome, so are $0w0$ and $1w1$.

Circumscription: Nothing else is a palindrome.

CFG's is a formal mechanism for definitions such as the one for L_{pal} .

1. $P \rightarrow \epsilon$
2. $P \rightarrow 0$
3. $P \rightarrow 1$
4. $P \rightarrow 0P0$
5. $P \rightarrow 1P1$

0 and 1 are *terminals*

P is a *variable* (or *nonterminal*, or *syntactic category*)

P is in this grammar also the *start symbol*.

1–5 are *productions* (or *rules*)

Formal definition of CFG's

A *context-free grammar* is a quadruple

$$G = (V, T, P, S)$$

where

V is a finite set of *variables*.

T is a finite set of *terminals*.

P is a finite set of *productions* of the form $A \rightarrow \alpha$, where A is a variable and $\alpha \in (V \cup T)^*$

S is a designated variable called the *start symbol*.

Example: $G_{pal} = (\{P\}, \{0, 1\}, A, P)$, where $A = \{P \rightarrow \epsilon, P \rightarrow 0, P \rightarrow 1, P \rightarrow 0P0, P \rightarrow 1P1\}$.

Sometimes we group productions with the same head, e.g. $A = \{P \rightarrow \epsilon|0|1|0P0|1P1\}$.

Example: Regular expressions over $\{0, 1\}$ can be defined by the grammar

$$G_{regex} = (\{E\}, \{0, 1\}, A, E)$$

where $A = \{ \text{+}, \text{.}, \text{\phi}, \text{\epsilon}, \text{*}, \text{(}, \text{)} \}$

$$\{E \rightarrow 0, E \rightarrow 1, E \rightarrow E.E, E \rightarrow E+E, E \rightarrow E^*, E \rightarrow (E)\}$$

$$E \rightarrow \epsilon, E \rightarrow \phi$$

Example: (simple) expressions in a typical prog lang. Operators are $+$ and $*$, and arguments are identifiers, i.e. strings in $L((a + b)(a + b + 0 + 1)^*)$

The expressions are defined by the grammar

$$G = (\{E, I\}, T, P, E)$$

where $T = \{+, *, (,), a, b, 0, 1\}$ and P is the following set of productions:

1. $E \rightarrow I$
2. $E \rightarrow E + E$
3. $E \rightarrow E * E$
4. $E \rightarrow (E)$
5. $I \rightarrow a$
6. $I \rightarrow b$
7. $I \rightarrow Ia$
8. $I \rightarrow Ib$
9. $I \rightarrow I0$
10. $I \rightarrow I1$

Derivations using grammars

- *Recursive inference*, using productions from body to head
- *Derivations*, using productions from head to body.

Example of recursive inference:

	String	Lang	Prod	String(s) used
(i)	a	I	5	-
(ii)	b	I	6	-
(iii)	$b0$	I	9	(ii)
(iv)	$b00$	I	9	(iii)
(v)	a	E	1	(i)
(vi)	$b00$	E	1	(iv)
(vii)	$a + b00$	E	2	(v), (vi)
(viii)	$(a + b00)$	E	4	(vii)
(ix)	$a * (a + b00)$	E	3	(v), (viii)

Let $G = (V, T, P, S)$ be a CFG, $A \in V$, $\{\alpha, \beta\} \subset (V \cup T)^*$, and $A \rightarrow \gamma \in P$.

Then we write

$$\alpha A \beta \xRightarrow[G]{} \alpha \gamma \beta$$

or, if G is understood

$$\alpha A \beta \Rightarrow \alpha \gamma \beta$$

and say that $\alpha A \beta$ *derives* $\alpha \gamma \beta$.

We define $\xRightarrow{*}$ to be the reflexive and transitive closure of \Rightarrow , IOW:

Basis: Let $\alpha \in (V \cup T)^*$. Then $\alpha \xRightarrow{*} \alpha$.

Induction: If $\alpha \xRightarrow{*} \beta$, and $\beta \Rightarrow \gamma$, then $\alpha \xRightarrow{*} \gamma$.

Example: Derivation of $a * (a + b00)$ from E in the grammar of slide 138:

$$\begin{aligned}
 E &\Rightarrow E * E \Rightarrow I * E \Rightarrow a * E \Rightarrow a * (E) \Rightarrow \\
 a*(E+E) &\Rightarrow a*(I+E) \Rightarrow a*(a+E) \Rightarrow a*(a+I) \Rightarrow \\
 a * (a + I0) &\Rightarrow a * (a + I00) \Rightarrow a * (a + b00)
 \end{aligned}$$

Note: At each step we might have several rules to choose from, e.g.

$$\begin{aligned}
 I * E &\Rightarrow a * E \Rightarrow a * (E), \text{ versus} \\
 I * E &\Rightarrow I * (E) \Rightarrow a * (E).
 \end{aligned}$$

Note2: Not all choices lead to successful derivations of a particular string, for instance

$$E \Rightarrow E + E$$

won't lead to a derivation of $a * (a + b00)$.

Leftmost and Rightmost Derivations

Leftmost derivation \Rightarrow : Always replace the leftmost variable by one of its rule-bodies.

Rightmost derivation \Rightarrow : Always replace the rightmost variable by one of its rule-bodies.

Leftmost: The derivation on the previous slide.

Rightmost:

$$E \xRightarrow{rm} E * E \xRightarrow{rm}$$

$$E*(E) \xRightarrow{rm} E*(E+E) \xRightarrow{rm} E*(E+I) \xRightarrow{rm} E*(E+I0)$$

$$\xRightarrow{rm} E*(E+I00) \xRightarrow{rm} E*(E+b00) \xRightarrow{rm} E*(I+b00)$$

$$\xRightarrow{rm} E*(a+b00) \xRightarrow{rm} I*(a+b00) \xRightarrow{rm} a*(a+b00)$$

We can conclude that $E \xRightarrow{*}{rm} a*(a+b00)$

The Language of a Grammar

If $G(V, T, P, S)$ is a CFG, then the *language of* G is

$$L(G) = \{w \in T^* : S \xrightarrow[G]{*} w\}$$

i.e. the set of strings over T^* derivable from the start symbol.

If G is a CFG, we call $L(G)$ a *context-free language*.

Example: $L(G_{pal})$ is a context-free language.

Theorem 5.7:

$$L(G_{pal}) = \{w \in \{0, 1\}^* : w = w^R\}$$

Proof: (\supseteq -direction.) Suppose $w = w^R$. We show by induction on $|w|$ that $w \in L(G_{pal})$

Basis: $|w| = 0$, or $|w| = 1$. Then w is $\epsilon, 0$, or 1 . Since $P \rightarrow \epsilon, P \rightarrow 0$, and $P \rightarrow 1$ are productions, we conclude that $P \xrightarrow[G]{*} w$ in all base cases.

Induction: Suppose $|w| \geq 2$. Since $w = w^R$, we have $w = 0x0$, or $w = 1x1$, and $x = x^R$.

If $w = 0x0$ we know from the IH that $P \xrightarrow{*} x$.
Then

$$P \Rightarrow 0P0 \xrightarrow{*} 0x0 = w$$

Thus $w \in L(G_{pal})$.

The case for $w = 1x1$ is similar.

(\subseteq -direction.) We assume that $w \in L(G_{pal})$ and must show that $w = w^R$.

Since $w \in L(G_{pal})$, we have $P \xRightarrow{*} w$.

We do an induction on the length of $\xRightarrow{*}$.

Basis: The derivation $P \xRightarrow{*} w$ is done in one step.

Then w must be ϵ , 0 , or 1 , all palindromes.

Induction: Let $n \geq 1$, and suppose the derivation takes $n + 1$ steps. Then we must have

$$w = 0x0 \xleftarrow{*} 0P0 \Leftarrow P$$

or

$$w = 1x1 \xleftarrow{*} 1P1 \Leftarrow P$$

where the second derivation is done in n steps.

By the IH x is a palindrome, and the inductive proof is complete.

Ex. Design CFGs for the following languages:

$L_1 = \{\text{balanced parentheses}\} = \{\epsilon, (), (()), ()(), ((())), ()(), ()(), \dots\}$ (the Dyck language)

$L_2 = \{0^m 1^n 2^p \mid m, n, p \geq 0, m+n = p\}$

$L_3 = \{w \mid w \in \{0,1\}^*, w \triangleleft w^R\}$

Sentential Forms

Let $G = (V, T, P, S)$ be a CFG, and $\alpha \in (V \cup T)^*$.

If

$$S \xRightarrow{*} \alpha$$

we say that α is a *sentential form*.

If $S \xRightarrow{lm} \alpha$ we say that α is a *left-sentential form*,
and if $S \xRightarrow{rm} \alpha$ we say that α is a *right-sentential form*

Note: $L(G)$ is those sentential forms that are in T^* .

Example: Take G from slide 138. Then $E * (I + E)$ is a sentential form since

$$E \Rightarrow E * E \Rightarrow E * (E) \Rightarrow E * (E + E) \Rightarrow E * (I + E)$$

This derivation is neither leftmost, nor rightmost

Example: $a * E$ is a left-sentential form, since

$$E \underset{lm}{\Rightarrow} E * E \underset{lm}{\Rightarrow} I * E \underset{lm}{\Rightarrow} a * E$$

Example: $E * (E + E)$ is a right-sentential form, since

$$E \underset{rm}{\Rightarrow} E * E \underset{rm}{\Rightarrow} E * (E) \underset{rm}{\Rightarrow} E * (E + E)$$