

Context-free Languages and Grammars

Prof. (Dr.) K.R. Chowdhary
Email: kr.chowdhary@iitj.ac.in

Formerly at department of Computer Science and Engineering
MBM Engineering College, Jodhpur

Saturday 10th December, 2016

- CFLs and CFGs are fundamentals to computer science, because they help in describing the structure of programming languages. All the HLL are in the category of CFL. Though natural languages (NLs) are not CFL, but their analysis is possible only when they are treated as CFLs.
- Context-free languages(CFL) are more powerful than regular languages. The context-free grammars (CFG) are generators of CFL. Regular languages are subset of CFL.
- CFG is finite specification of rules to generate infinite context-free language.
- To generate all the strings of regex $a^*(b^* + c^*)$, we follow the steps:
 - 1 write character a zero or more times
 - 2 arbitrarily choose b or c and write it arbitrary times
 - 3 write c .

Generating Language strings

- Let $L = \mathcal{L}(a^*(b^* + c^*))$ is language corresponding to regex. We can generate all the strings by **substitution** rules:

a. $S \rightarrow AMb$, f. $B \rightarrow \varepsilon$

b. $A \rightarrow \varepsilon$, g. $B \rightarrow bB$

c. $A \rightarrow aA$, h. $C \rightarrow \varepsilon$

d. $M \rightarrow B$, i. $C \rightarrow \varepsilon$

e. $M \rightarrow C$, j. $C \rightarrow cC$

Consider generating $w = aaccb$

using production rules.

$$\begin{aligned} S &\Rightarrow AMb \text{ ;by rule a} \\ &\Rightarrow aAMb \text{ ;by rule c} \\ &\Rightarrow aaAMb \text{ ;by rule c} \\ &\Rightarrow aaMb \text{ ;by rule b} \\ &\Rightarrow aaCb \text{ ;by rule e} \\ &\Rightarrow aaccCb \text{ ;by rule j} \\ &\Rightarrow aaccb \text{ ;by rule i} \end{aligned}$$

Therefore, $aaccb \in L$. Symbol “ \rightarrow ” stand for “**can be substituted by**”, and “ \Rightarrow ” stand for “**derives.**” Strings like $aacCb$ or AMb , during the derivation are called sentential form.

Generating language strings

Let us try to generate the strings of language $L = \{a^n b^n | n \geq 0\}$. The rules this time are: $S \rightarrow aSb$, $S \rightarrow \epsilon$. Consider deriving $w = aaabbb$.

$$\begin{aligned} S &\Rightarrow aSb \\ &\Rightarrow aaSbb \\ &\Rightarrow aaaSbbb \\ &\Rightarrow aaabbb \end{aligned}$$

Therefore, $S \Rightarrow^* aaabbb$. The generator of these languages, is CFG $G = \{V, \Sigma, S, P\}$, where:

- V is finite set of variables

symbols, appearing in the process of derivation

- Σ is set of terminal symbols (appearing in the final generated sentence),
 $V \cap \Sigma = \phi$,
- S is start symbol,
- P is set of production/substitution rules of the form $A \rightarrow \alpha$, where $A \in V$ and $\alpha \in (V \cup \Sigma)^*$.
- Symbols in upper case in the beginning of English alphabets are variables (non-terminal symbols), i.e. A, B, C, D, E .

Definition

(Context-free grammar:) A CFG is regular grammar if productions are like: $A \rightarrow a, A \rightarrow aB, A \rightarrow \varepsilon$, where, $a \in \Sigma$, and $A, B \in V$. For the regular expression $a^*(b^* + c^*)b$, a CFG was used in the previous slides to generate the regular language.

Definition

(Derivation:) A derivation $\alpha_1 \Rightarrow_G \alpha_2 \Rightarrow_G \dots \Rightarrow_G \alpha_n$, can be written in short as $\alpha_1 \Rightarrow_G^* \alpha_n$. In a derivation: $\beta A \gamma \Rightarrow_G \beta \alpha \gamma$, the symbol A can be always substituted by α , if there is production like $A \rightarrow \alpha$, irrespective of presence of substrings β and γ around the non-terminal symbol A . Languages having this property are called **context-free**. The substrings β and γ are called context of variable A . Relation \Rightarrow is *reflexive, anti-symmetric, and transitive* (a partial ordering) relation.

Definition

Language acceptability: $L = L(G) = \{w \in \Sigma^* \mid S \Rightarrow_G^* w\}$. Two are equal if they generate the same language.

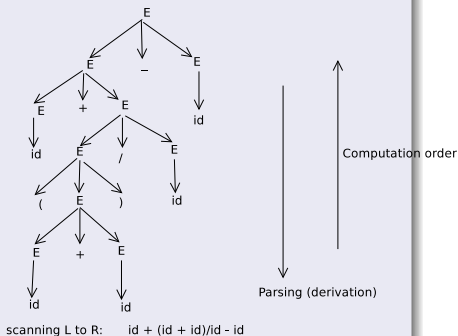
Derivations

Example

Given grammar $G = \{V, \Sigma, S, P\}$, $\Sigma = \{+, -, *, /, (,), id\}$, $V = \{E\}$, $S = E$, and $P = \{E \rightarrow E + E \mid E - E \mid E / E \mid E * E \mid (E) \mid id\}$, find out the derivation and derivation tree for $id * (id + id) - id$.

Solution

The generating process:

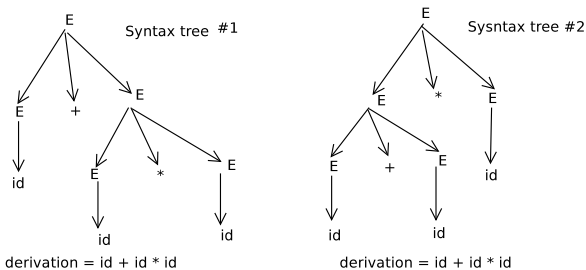
$$\begin{aligned} E &\Rightarrow E - E \Rightarrow E * E - E \\ &\Rightarrow id * E / E - E \Rightarrow id * (E) / E - E \\ &\Rightarrow id * (E + E) / E - E \Rightarrow id * \\ &\quad (id + E) / E - E \\ &\Rightarrow id * (id + id) / E - E \\ &\Rightarrow id * (id + id) / id - E \\ &\Rightarrow id * (id + id) / id - id \end{aligned}$$


Derivations and ambiguity

- Compilers derive a given expression, if it succeed, the expression is syntactically correct, else not.
 - In a derivation, e.g., $S \Rightarrow ABC$, we may start by first replacing left side variables (left hand derivation) or the right side variable first (right hand derivation). In both, the end result is the same.
 - If a language $L = L(G)$ can be derived using two or more different derivation trees, the grammar G is **ambiguous grammar**.
- If G has maximum n number of derivation trees, then its degree of ambiguity n .
 - It is recursively unsolvable, to find out if an arbitrary grammar is ambiguous. Hence, there does not exist an algorithm to find out if a given grammar is ambiguous.
 - A grammar is un-ambiguous if every $w \in L(G)$ has a unique parse-tree. A grammar is called **reduced**, if every non-terminal appears in some derivation.

Ambiguity

- Show that grammar for $id + id * id$ is ambiguous.



- The two derivation trees have different semantics: first calculates $id + (id * id)$ while the second does $(id + id) * id$, hence the grammar is ambiguous.

Removing ambiguity

The general case of detection of ambiguity in a grammar is unsolvable. However, if it is found that the grammar is ambiguous, it can be made unambiguous by adding few more non-terminals in the grammar.

Example

Given $\Sigma = \{ (,), +, *, id \}$, $P = \{ E \rightarrow E + E \mid E * E \mid (E) \mid id \}$, which is an ambiguous grammar, find out its equivalent unambiguous grammar.

Solution

Let $V = \{ E, T, F \}$, and
 $P = \{ E \rightarrow T, T \rightarrow F, F \rightarrow id, E \rightarrow E + T, T \rightarrow T * F, F \rightarrow (E) \}$. Note that, you can derive a string in one way only.

$$\begin{aligned} E &\Rightarrow E + T \\ &\Rightarrow T + T \\ &\Rightarrow F + T \Rightarrow id + T \\ &\Rightarrow id + T * F \Rightarrow id + F * F \\ &\Rightarrow id + id * F \Rightarrow id + id * id \end{aligned}$$