

بسم الله الرحمن الرحيم

دانشگاه صنعتی اصفهان - دانشکده مهندسی برق و کامپیوتر
(نیم سال تحصیلی ۴۰۲۱)

نظریه زبانها و ماشینها

حسین فلسفین

Leftmost and Rightmost Derivations

A derivation of a string w in a grammar G is a leftmost derivation if at every step the leftmost remaining variable is the one replaced. In other words, a leftmost derivation is a derivation in which the variable replaced at each step is the leftmost one. A rightmost derivation is defined analogously.

Example: In the grammar G with productions

$$S \rightarrow a | S + S | S * S | (S),$$

the string $a + a * a$ has the following two leftmost derivations:

$$S \Rightarrow S + S \Rightarrow a + S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * a$$

$$S \Rightarrow S * S \Rightarrow S + S * S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * a$$

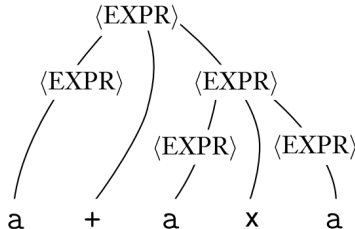
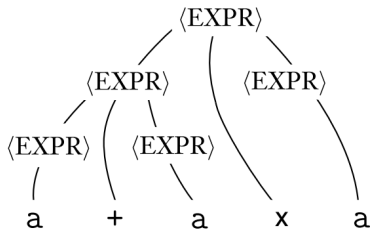
Ambiguity

Sometimes a grammar can generate the same string in several different ways. Such a string will have **several different parse trees** and thus several different meanings. This result may be **undesirable** for certain applications, such as programming languages, where a program should have a unique interpretation.

Example: consider grammar:

$$\langle \text{EXPR} \rangle \rightarrow \langle \text{EXPR} \rangle + \langle \text{EXPR} \rangle | \langle \text{EXPR} \rangle \times \langle \text{EXPR} \rangle | (\langle \text{EXPR} \rangle) | a$$

This grammar generates the string $a + a \times a$ ambiguously. The following figure shows the two different parse trees.



*In contrast, the following grammar generates exactly the same language, but every generated string has a unique parse tree. Hence it is **unambiguous**, whereas the previous one was **ambiguous**.*

$$\begin{aligned}\langle \text{EXPR} \rangle &\rightarrow \langle \text{EXPR} \rangle + \langle \text{TERM} \rangle \mid \langle \text{TERM} \rangle \\ \langle \text{TERM} \rangle &\rightarrow \langle \text{TERM} \rangle \times \langle \text{FACTOR} \rangle \mid \langle \text{FACTOR} \rangle \\ \langle \text{FACTOR} \rangle &\rightarrow (\langle \text{EXPR} \rangle) \mid a\end{aligned}$$

A context-free grammar G is said to be **ambiguous** if there exists some $w \in L(G)$ that has at least two distinct **derivation trees**.

- * When we say that a grammar generates a string ambiguously, we mean that the string has two different parse trees, not two different derivations. **Two derivations may differ merely in the order in which they replace variables yet not in their overall structure.**
- * Derivation trees show which productions are used in obtaining a sentence, but do not give the order of their application. A parse tree represents the rules applied, **but not the order of their application.**
- * A grammar is **unambiguous** if and only if every word generated has exactly one parse tree.

تذکر: درایویشن‌هایی که درخت تجزیه آنها یکسان است، ماهیتاً تمایز عمده‌ای با هم ندارند و صرفاً ترتیب بکارگیری رول‌ها در آنها متفاوت است. در مقابل، درایویشن‌هایی که درخت‌های تجزیه آنها متفاوت است، ماهیتاً متفاوت هستند. مثلاً، گرامر

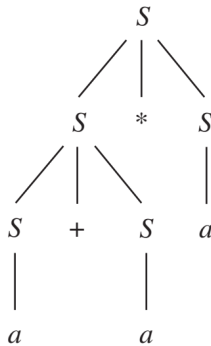
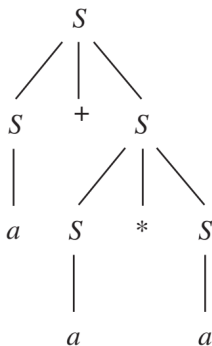
$$S \rightarrow a|S + S|S * S|(S)$$

را در نظر بگیرید.

در این گرامر، دو درایویشن زیر برای رشته $a + a \times a$ ، ماهیتاً متفاوت هستند چون درخت‌های تجزیه نظیر آن دو متفاوت است:

$$S \Rightarrow S + S \Rightarrow a + S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * a$$

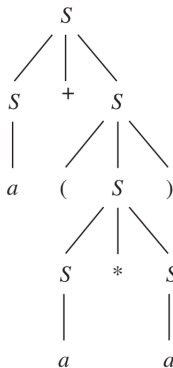
$$S \Rightarrow S * S \Rightarrow S + S * S \Rightarrow a + S * S \Rightarrow a + a * S \Rightarrow a + a * a$$



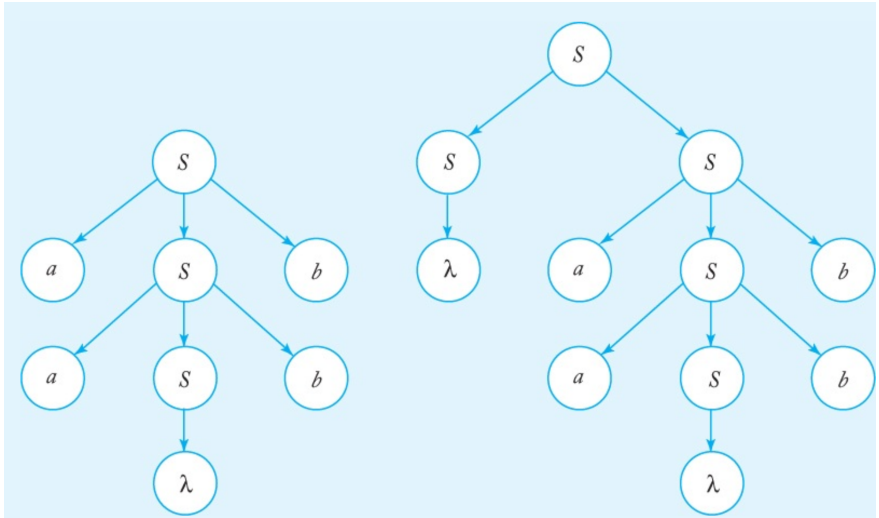
اما دو درایویشن زیر برای رشته $a + (a \times a)$ ماهیتاً متمایز نیستند چون درخت تجزیه نظیر آنها یکسان است. تمایز این دو صرفاً در ترتیب بکارگیری رول‌ها است:

$$S \Rightarrow S+S \Rightarrow a+S \Rightarrow a+(S) \Rightarrow a+(S*S) \Rightarrow a+(a*S) \Rightarrow a+(a*a)$$

$$S \Rightarrow S+S \Rightarrow S+(S) \Rightarrow S+(S*S) \Rightarrow S+(a*S) \Rightarrow S+(a*a) \Rightarrow a+(a*a)$$



Example: The grammar $S \rightarrow aSb | SS | \varepsilon$ is ambiguous. The sentence $aabb$ has the two derivation trees:



مفهوم ابهام با بهره‌گیری از دو مفهوم LMD و RMD نیز قابل بیان است:

If G is a context-free grammar, then for every $x \in L(G)$, these three statements **are equivalent**:

1. x has more than one derivation tree.
2. x has more than one leftmost derivation.
3. x has more than one rightmost derivation.

Alternatively, ambiguity implies the existence of two or more LMDs (RMDs). In fact, a string w is derived ambiguously in CFG G if it has two or more different LMDs (RMDs). Grammar G is ambiguous if it generates some string ambiguously. A grammar G is said to be unambiguous if every word $w \in L(G)$ has exactly one LMD (RMD) and ambiguous otherwise.

مفهوم ذاتاً مبهم بودن که به یک زبان اطلاق می‌شود:

If L is a context-free language for which there exists an unambiguous grammar, then L is said to be unambiguous. If **every** grammar that generates L is ambiguous, then the language is called **inherently ambiguous**. Indeed, sometimes when we have an ambiguous grammar we can find an unambiguous grammar that generates the same language. Some context-free languages, however, can be generated **only** by ambiguous grammars. Such languages are called **inherently ambiguous**.

It can be shown that, the language $\{a^i b^j c^k \mid i = j \text{ or } j = k\}$ is inherently ambiguous.

Example: The language $L = \{a^n b^n c^m\} \cup \{a^n b^m c^m\}$, with n and m nonnegative, is an inherently ambiguous context-free language. That L is context-free is easy to show. Notice that $L = L_1 \cup L_2$, where L_1 is generated by

$$S_1 \rightarrow S_1 c | A, \quad A \rightarrow a A b | \varepsilon$$

and L_2 is given by an analogous grammar with start symbol S_2 and productions

$$S_2 \rightarrow a S_2 | B, \quad B \rightarrow b B c | \varepsilon$$

Then L is generated by the combination of these two grammars with the additional production $S \rightarrow S_1 | S_2$. The grammar is ambiguous since the string $a^n b^n c^n$ has two distinct derivations, one starting with $S \Rightarrow S_1$, the other with $S \Rightarrow S_2$. It does not, of course, follow from this that L is inherently ambiguous as there might exist some other unambiguous grammars for it. A rigorous argument is quite technical.

Chomsky Normal Form

When working with context-free grammars, it is often convenient to have them in simplified form. One of the simplest and most useful forms is called the Chomsky normal form. In general, a context-free grammar can be transformed to a grammar in Chomsky normal form without changing the language to be generated.

تعریف فرم نرمال چامسکی

A context-free grammar is in Chomsky normal form if every rule is of the form

$$A \rightarrow BC$$

$$A \rightarrow a$$

where a is any terminal and A , B , and C are any variables—except that B and C may not be the start variable. In addition, we permit the rule $S \rightarrow \varepsilon$, where S is the start variable.

The ε -rule $S \rightarrow \varepsilon$ is permitted only for the start symbol S .

Theorem: Any context-free language is generated by a context-free grammar in Chomsky normal form.

Proof Idea: We can convert any grammar G into Chomsky normal form. The conversion has several stages wherein rules that violate the conditions are replaced with equivalent ones that are satisfactory.

👉 First, we add a new start variable.

👉 Then, we eliminate all ε -rules of the form $A \rightarrow \varepsilon$.

👉 We also eliminate all unit (renaming) rules of the form $A \rightarrow B$.

👉 In both cases we patch up the grammar to be sure that it still generates the same language.

👉 Finally, we convert the remaining rules into the proper form. A rule $A \rightarrow u_1 u_2 \cdots u_k$ in which the length of the right-hand side k is more than or equal to 3 is called a long rule.

پس روند در مجموع به صورت ذیل است:

We transform G to an equivalent context-free grammar in Chomsky normal form by **deleting the ϵ -rules, the renaming rules, and finally the long rules** from the rules of G . We shall explain the idea for deleting the three types of rules.

Proof:

☞ **First**, we add a new start variable S_0 and the rule $S_0 \rightarrow S$, where S was the original start variable. This change guarantees that the start variable doesn't occur on the right-hand side of a rule. Observe that this change does not change the language to be generated and that the start symbol S_0 never appears on the right-hand side of any rule.

☞ **Second**, we take care of all ε -rules. We remove an ε -rule $A \rightarrow \varepsilon$, where A is not the start variable. Then for each occurrence of an A on the right-hand side of a rule, we add a new rule with that occurrence deleted. In other words, if $R \rightarrow uAv$ is a rule in which u and v are strings of variables and terminals, we add rule $R \rightarrow uv$. We do so for each occurrence of an A , so the rule $R \rightarrow uAvAw$ causes us to add $R \rightarrow uvAw$, $R \rightarrow uAvw$, and $R \rightarrow uvw$. If we have the rule $R \rightarrow A$, we add $R \rightarrow \varepsilon$ unless we had previously removed the rule $R \rightarrow \varepsilon$. We repeat these steps until we eliminate

all ε -rules not involving the start variable.

☞ **Third**, we handle all unit rules (deletion of renaming rules). We remove a unit rule $A \rightarrow B$. Then, whenever a rule $B \rightarrow u$ appears, we add the rule $A \rightarrow u$ unless this was a unit rule previously removed. As before, u is a string of variables and terminals. We repeat these steps until we eliminate all unit rules. In fact, for each renaming rule $A \rightarrow B$, the deletion of $A \rightarrow B$ is applied for all the rules that take the form of $B \rightarrow u$, where $u \in (V \cup \Sigma)^*$. By the deletion, we eliminate the rule $A \rightarrow B$ and add all the rules that are expressed as $A \rightarrow u$ provided that $A \rightarrow u$ was not deleted so far. For each renaming rule this type of deletion and addition will be applied concerning all the remaining rules, including ones that are added while applying other deletions.

☞ **Finally**, we convert all remaining rules into the proper form. We replace each rule $A \rightarrow u_1 u_2 \cdots u_k$, where $k \geq 3$ and each u_i is a variable or terminal symbol, with the rules $A \rightarrow u_1 A_1$, $A_1 \rightarrow u_2 A_2$,

$A_2 \rightarrow u_3 A_3, \dots, \text{ and } A_{k-2} \rightarrow u_{k-1} u_k.$ The A_i 's are new variables. We replace any terminal u_i in the preceding rule(s) with the new variable U_i and add the rule $U_i \rightarrow u_i$.

نکات مهمی درباره حذف ϵ -rule ها که باید مجدداً متذکر شوم:

- ☞ Since in Chomsky normal form $S \rightarrow \epsilon$ is permitted for the start symbol S , the ϵ -rules we must delete are the ones with variables on the left-hand side other than the start symbol.
- ☞ Once the deletion of $A \rightarrow \epsilon$ is applied for all the rules in the grammar, the ϵ -rule $A \rightarrow \epsilon$ will never be taken up again in the later deletion process.
- ☞ After the deletion of all the ϵ rules for all the rules is done, the only ϵ rule that may possibly be left is $S_0 \rightarrow \epsilon$.

Example: The series of grammars presented illustrates the steps in the conversion. Rules shown in bold have just been added. Rules shown in gray have just been removed.

1. The original CFG G_6 is shown on the left. The result of applying the first step to make a new start variable appears on the right.

$S \rightarrow ASA \mid aB$ $A \rightarrow B \mid S$ $B \rightarrow b \mid \epsilon$	$S_0 \rightarrow S$ $S \rightarrow ASA \mid aB$ $A \rightarrow B \mid S$ $B \rightarrow b \mid \epsilon$
--	--

2. Remove ϵ -rules $B \rightarrow \epsilon$, shown on the left, and $A \rightarrow \epsilon$, shown on the right.

$S_0 \rightarrow S$ $S \rightarrow ASA \mid aB \mid \mathbf{a}$ $A \rightarrow B \mid S \mid \epsilon$ $B \rightarrow b \mid \epsilon$	$S_0 \rightarrow S$ $S \rightarrow ASA \mid aB \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS} \mid S$ $A \rightarrow B \mid S \mid \epsilon$ $B \rightarrow b$
--	---

ادامه مثال:

3a. Remove unit rules $S \rightarrow S$, shown on the left, and $S_0 \rightarrow S$, shown on the right.

$S_0 \rightarrow S$ $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS \mid S$ $A \rightarrow B \mid S$ $B \rightarrow b$	$S_0 \rightarrow \mathbf{S} \mid \mathbf{ASA} \mid \mathbf{aB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$ $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$ $A \rightarrow B \mid S$ $B \rightarrow b$
--	---

3b. Remove unit rules $A \rightarrow B$ and $A \rightarrow S$.

$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$ $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$ $A \rightarrow \mathbf{B} \mid \mathbf{S} \mid \mathbf{b}$ $B \rightarrow b$	$S_0 \rightarrow ASA \mid aB \mid a \mid SA \mid AS$ $S \rightarrow ASA \mid aB \mid a \mid SA \mid AS$ $A \rightarrow \mathbf{S} \mid \mathbf{b} \mid \mathbf{ASA} \mid \mathbf{aB} \mid \mathbf{a} \mid \mathbf{SA} \mid \mathbf{AS}$ $B \rightarrow b$
--	---

ادامه مثال:

4. Convert the remaining rules into the proper form by adding additional variables and rules. The final grammar in Chomsky normal form is equivalent to G_6 . (Actually the procedure given in Theorem 2.9 produces several variables U_i and several rules $U_i \rightarrow a$. We simplified the resulting grammar by using a single variable U and rule $U \rightarrow a$.)

$$\begin{aligned} S_0 &\rightarrow AA_1 \mid UB \mid a \mid SA \mid AS \\ S &\rightarrow AA_1 \mid UB \mid a \mid SA \mid AS \\ A &\rightarrow b \mid AA_1 \mid UB \mid a \mid SA \mid AS \\ A_1 &\rightarrow SA \\ U &\rightarrow a \\ B &\rightarrow b \end{aligned}$$



Example: We transform the grammar $S \rightarrow aS|aSbS|\varepsilon$ to an equivalent context-free grammar which is in Chomsky normal form in the way described above.

First, introducing the start symbol S_0 , we modify the grammar as follows:

$$S_0 \rightarrow S, S \rightarrow aS|aSbS|\varepsilon.$$

Next, we delete $S \rightarrow \varepsilon$ by applying what is described above to the three rules: By applying to $S_0 \rightarrow S$ we have $S_0 \rightarrow S|\varepsilon$. By applying to $S \rightarrow aS$ we have $S \rightarrow aS|a$. By applying to $S \rightarrow aSbS$ we have $S \rightarrow aSbS|abS|aSb|ab$. Putting them together, we have

$$S_0 \rightarrow S|\varepsilon, S \rightarrow aS|a|aSbS|abS|aSb|ab.$$

Next, to delete $S_0 \rightarrow S$, we add $S_0 \rightarrow w$ for each rule of the form $S \rightarrow w$, thereby obtaining rules as follows:

$$S_0 \rightarrow aS|a|aSbS|abS|aSb|ab|\varepsilon, S \rightarrow aS|a|aSbS|abS|aSb|ab.$$

Finally, we delete long rules. Deleting $S \rightarrow aSbS$, we add instead

$$S \rightarrow X_a A_2, A_2 \rightarrow S A_3, A_3 \rightarrow X_b S, X_a \rightarrow a, \text{ and } X_b \rightarrow b.$$

In this way we have the equivalent rules in Chomsky normal form as follows:

$$\begin{aligned} S_0 &\rightarrow X_a S | a | X_a A_2 | X_a B_2 | X_a C_2 | X_a X_b | \varepsilon, \\ S &\rightarrow X_a S | a | X_a A_2 | X_a B_2 | X_a C_2 | X_a X_b, \\ A_2 &\rightarrow S A_3, \\ A_3 &\rightarrow X_b S, \\ B_2 &\rightarrow X_b S, \\ C_2 &\rightarrow S X_b, \\ X_a &\rightarrow a, \\ X_b &\rightarrow b. \end{aligned}$$

تمرین:

Transform the grammar

$$S \rightarrow aSaaA|A, \quad A \rightarrow abA|bb$$

into Chomsky normal form.

Greibach Normal Form

Another useful grammatical form is the Greibach normal form. **Here we put restrictions not on the length of the right sides of a production, but on the positions in which terminals and variables can appear.** Arguments justifying Greibach normal form are a little complicated and not very transparent. Similarly, constructing a grammar in Greibach normal form equivalent to a given context-free grammar is tedious. We therefore deal with this matter very briefly. Nevertheless, the Greibach normal form has several uses, both practical and theoretical.

Any CFG can be written in a special form called Greibach normal form. (For every CFG G , there exists an equivalent grammar \hat{G} in Greibach normal form.) **Each production has one of the following forms: $A \rightarrow aB_1B_2 \cdots B_n$, $A \rightarrow a$, $S \rightarrow \varepsilon$, where S is the start symbol, a is a terminal, and each B_k is a variable not equal to S . ε can occur only with the production $S \rightarrow \varepsilon$.**