

15-453 Formal Languages, Automata, and Computation

Some Notes on Regular Grammars

Frank Pfenning

Lecture 9
February 7, 2000

In this notes we describe restrictions to context-free grammars which ensure that the generated languages are regular. We first discuss *strictly right-linear grammars* and show that they correspond directly to non-deterministic finite automata (NFAs). We then generalize to *right-linear grammars* which does not change the set of accepted languages, but allows more concise specifications. The dual *left-linear grammars* also generate the same languages. A grammar is *regular* if it is either right-linear or left-linear.

1 Strictly Right-Linear Grammars

We first review the definition of a *context-free grammar*.

Definition 1 (Context-Free Grammar)

A context-free grammar G is specified by a tuple (V, Σ, R, S) where

V	is a finite set of variables (also called non-terminals),
Σ	is a finite alphabet,
$R \subseteq V \times (V \cup \Sigma)^*$	is the set of rules,
$S \in V$	is the start variable.

We write rules (A, w) as $A \rightarrow w$. A derivation step has the form $uAv \Rightarrow uwv$ where $A \in V$ and $A \rightarrow w \in R$. A derivation has the form

$$u_1 \Rightarrow u_2 \Rightarrow \dots \Rightarrow u_n$$

which we abbreviate $u_1 \xrightarrow{*} u_n$. Finally, we define the language generated by G

$$\mathcal{L}(G) = \{w \in \Sigma^* \mid S \xrightarrow{*} w\}.$$

For a strictly right-linear grammar, we severely restrict the form of the rules.

Definition 2 (Strictly Right-Linear Grammar) A strictly right-linear grammar is a context-free grammar G where each rule has one of the following forms:

$$\begin{aligned} A \rightarrow xB & \text{ for } x \in \Sigma_\epsilon \\ A \rightarrow \epsilon & \end{aligned}$$

Any derivation of a word w from a start symbol S in a strictly right-linear grammar has the form

$$S \Rightarrow x_1 A_1 \Rightarrow x_1 x_2 A_2 \Rightarrow \dots \Rightarrow x_1 \dots x_n A_n \Rightarrow x_1 \dots x_n$$

where $w = x_1 \dots x_n$. Note that some of the x_i may be ϵ . The form of these derivations suggests the connection to non-deterministic finite automata.

2 Strictly Right-Linear Grammars Generate Regular Languages

Lemma 1 *If G is a strictly right-linear grammar, then $\mathcal{L}(G)$ is regular.*

Proof: Given a strictly right-linear grammar $G = (V, \Sigma, R, S)$, we construct a non-deterministic finite automaton $N = (V, \Sigma, \delta, S, F)$ such that $\mathcal{L}(G) = \mathcal{L}(N)$. The set of states is simply the set of non-terminals of G , with the start state corresponding to the start variable. We further define

$$\begin{aligned}\delta(A, x) &= \{B \mid A \rightarrow xB \in R\}, \\ F &= \{C \mid C \rightarrow \epsilon \in R\}.\end{aligned}$$

It remains to show that $\mathcal{L}(G) = \mathcal{L}(N)$. As usual, we do this in two parts, $\mathcal{L}(G) \subseteq \mathcal{L}(N)$ and $\mathcal{L}(N) \subseteq \mathcal{L}(G)$. We write $A \xrightarrow{w} N B$ for computations of the NFA N .

1. $\mathcal{L}(G) \subseteq \mathcal{L}(N)$. We show by induction on the structure of the derivation:

If $A \xrightarrow{*} w$ then either

- (a) $w = uB$ and $A \xrightarrow{u} N B$ for some $u \in \Sigma^*$, or
- (b) $w = u$ and $A \xrightarrow{u} N C$ for some $u \in \Sigma^*$ and final state C .

Base Case: $A \xrightarrow{*} A$. Then $A \xrightarrow{\epsilon} N A$ by definition of $\xrightarrow{*}$.

Induction Step, Case 1: $A \Rightarrow x A_1 \xrightarrow{*} x w_1$ for $x \in \Sigma_\epsilon$ and $w = x w_1$. By induction hypothesis on $A_1 \xrightarrow{*} w_1$, we have two subcases to consider:

- (a) $A_1 \xrightarrow{*} vB$ and $A_1 \xrightarrow{v} N B$. Then $u = xv$ satisfies the claim, since $A \Rightarrow xvB$ and $A \xrightarrow{x} N A_1 \xrightarrow{v} N B$.
- (b) $A_1 \xrightarrow{*} v$ and $A_1 \xrightarrow{u} N C$ for some final state C . Then $u = xv$ satisfies the claim, since $A \Rightarrow xv$ and $A \xrightarrow{x} N A_1 \xrightarrow{v} N C$ where C is a final state.

Induction Step, Case 2: $A \Rightarrow \epsilon$. Then $A \in F$ and $A \xrightarrow{\epsilon} N A$.

2. $\mathcal{L}(N) \subseteq \mathcal{L}(G)$. We show by induction on the structure of the computation:

If $A \xrightarrow{w} N B$ then $A \xrightarrow{*} wB$.

From this the claim follows by using S for A and a final state C for B .

Base Case: $A \xrightarrow{\epsilon} N A$. Then $A \xrightarrow{*} A$ by definition of $\xrightarrow{*}$.

Induction Step: $A \xrightarrow{x} N A_1 \xrightarrow{u} N B$ where $w = xu$ and $x \in \Sigma_\epsilon$. By induction hypothesis, $A_1 \xrightarrow{*} uB$ so $A \Rightarrow x A_1 \xrightarrow{*} xuB$.

□

Lemma 2 *If L is a regular language, then there is a strictly right-linear grammar G such that $\mathcal{L}(G) = L$.*

Proof: If L is regular, then there is an NFA $N = (Q, \Sigma, \delta, q_0, F)$ recognizing L . We construct a strictly right-linear grammar $G = (Q, \Sigma, R, q_0)$ where

$$R = \{q \rightarrow xr \mid r \in \delta(q, x)\} \cup \{q \rightarrow \epsilon \mid q \in F\}.$$

We leave the proof that $\mathcal{L}(N) = \mathcal{L}(G)$ to the reader. It follows by straightforward inductions in both directions. \square

Theorem 1 *A language L is regular if and only if there is a strictly right-linear grammar G such that $\mathcal{L}(G) = L$.*

Proof: By the two preceding lemmas. \square

3 Right-Linear Grammars

Right-linear grammars allow more general right-hand sides for rules, but they can easily be translated to strictly right-linear grammars.

Definition 3 (Right-Linear Grammar)

A right-linear grammar is a context-free grammar $G = (V, \Sigma, R, S)$ where each rule in R has one of the following forms:

$$\begin{aligned} A &\rightarrow wB & \text{for } w \in \Sigma_\epsilon^* \\ A &\rightarrow w & \text{for } w \in \Sigma_\epsilon^* \end{aligned}$$

Clearly, any strictly right-linear grammar is right-linear. In addition, strictly right-linear grammars can easily simulate right-linear grammars.

Lemma 3 *For any right-linear grammar G there exists a strictly right-linear grammar H such that $\mathcal{L}(G) = \mathcal{L}(H)$.*

Proof: Let $G = (V, \Sigma, R, S)$ be a right-linear grammar. We construct a strictly right-linear grammar by transforming each rule

$$A \rightarrow a_1 \dots a_n B$$

where $n \geq 2$ into the set of rules

$$\begin{aligned} A &\rightarrow a_1 A_1 \\ A_1 &\rightarrow a_2 A_2 \\ &\dots \\ A_n &\rightarrow B \end{aligned}$$

where A_1, \dots, A_n are new non-terminals. Similarly, each rule

$$A \rightarrow a_1 \dots a_n$$

for $n \geq 1$ is translated into

$$\begin{aligned} A &\rightarrow a_1 A_1 \\ A_1 &\rightarrow a_2 A_2 \\ &\dots \\ A_n &\rightarrow \epsilon \end{aligned}$$

where A_1, \dots, A_n are new non-terminals.

Then $H = (V', \Sigma, R', S)$, where V' contains V and all new non-terminals introduced in the transformation above, and R' contains the rules of R satisfying the strict right-linearity condition and all rules resulting from the above transformation.

It is now easy to see that $\mathcal{L}(G) = \mathcal{L}(H)$. We leave the details of an inductive proof to the reader. \square

Theorem 2 *A language L is regular if and only if there is a right-linear grammar G with $\mathcal{L}(G) = L$.*

Proof: By the preceding lemma and Theorem 1. \square

4 Regular Grammars

Regular grammars are grammars that are either right-linear or left-linear as defined below.

Definition 4 (Left-Linear Grammars)

A left-linear grammar is a context-free grammar $G = (V, \Sigma, R, S)$ where each rule in R has one of the following forms:

$$\begin{aligned} A &\rightarrow Bw & \text{for } w \in \Sigma_\epsilon^* \\ A &\rightarrow w & \text{for } w \in \Sigma_\epsilon^* \end{aligned}$$

It is easy to show that left-linear grammars also define regular languages. This leads to the theorem that a language L is regular if and only if it is generated by a regular grammar.

Note, however, that if we are allowed to mix left-linear and right-linear rules in a single grammar, the result will not necessarily generate a regular language. For example, a grammar with the rules

$$\begin{aligned} S &\rightarrow 0A \\ S &\rightarrow 1B \\ S &\rightarrow \epsilon \\ A &\rightarrow S0 \\ B &\rightarrow S1 \end{aligned}$$

generates the language $L = \{ww^R \mid w \in \{0, 1\}^*\}$ which is not regular.