



Regular grammar

In theoretical computer science and formal language theory, a **regular grammar** is a grammar that is *right-regular* or *left-regular*. While their exact definition varies from textbook to textbook, they all require that

- all production rules have at most one non-terminal symbol;
- that symbol is either always at the end or always at the start of the rule's right-hand side.

Every regular grammar describes a regular language.

Strictly regular grammars

A **right-regular grammar** (also called right-linear grammar) is a formal grammar (N, Σ, P, S) in which all production rules in P are of one of the following forms:

1. $A \rightarrow a$
2. $A \rightarrow aB$
3. $A \rightarrow \varepsilon$

where $A, B, S \in N$ are non-terminal symbols, $a \in \Sigma$ is a terminal symbol, and ε denotes the empty string, i.e. the string of length 0. S is called the start symbol.

In a **left-regular grammar**, (also called left-linear grammar), all rules obey the forms

1. $A \rightarrow a$
2. $A \rightarrow Ba$
3. $A \rightarrow \varepsilon$

The language described by a given grammar is the set of all strings that contain only terminal symbols and can be derived from the start symbol by repeated application of production rules. Two grammars are called weakly equivalent if they describe the same language.

Rules of both kinds must not be mixed; for example, the grammar with rule set $\{ S \rightarrow aT, T \rightarrow Sb, S \rightarrow \varepsilon \}$ is not regular, and describes the language $\{a^i b^i : i \in \mathbb{N}\}$, which is not regular, either.

Some textbooks and articles disallow empty production rules, and assume that the empty string is not present in languages.

Extended regular grammars

An *extended right-regular grammar* is one in which all rules obey one of

1. $A \rightarrow w$, where A is a non-terminal in N and w is in a (possibly empty) string of terminals Σ^*
2. $A \rightarrow wB$, where A and B are in N and w is in Σ^* .

Some authors call this type of grammar a *right-regular grammar* (or *right-linear grammar*)^[1] and the type above a *strictly right-regular grammar* (or *strictly right-linear grammar*).^[2]

An *extended left-regular grammar* is one in which all rules obey one of

1. $A \rightarrow w$, where A is a non-terminal in N and w is in Σ^*
2. $A \rightarrow Bw$, where A and B are in N and w is in Σ^* .

Examples

An example of a right-regular grammar G with $N = \{S, A\}$, $\Sigma = \{a, b, c\}$, P consists of the following rules

$S \rightarrow aS$
 $S \rightarrow bA$
 $A \rightarrow \varepsilon$
 $A \rightarrow cA$

and S is the start symbol. This grammar describes the same language as the regular expression a^*bc^* , viz. the set of all strings consisting of arbitrarily many " a "s, followed by a single " b ", followed by arbitrarily many " c "s.

A somewhat longer but more explicit extended right-regular grammar G for the same regular expression is given by $N = \{S, A, B, C\}$, $\Sigma = \{a, b, c\}$, where P consists of the following rules:

$S \rightarrow A$
 $A \rightarrow aA$
 $A \rightarrow B$
 $B \rightarrow bC$
 $C \rightarrow \varepsilon$
 $C \rightarrow cC$

...where each uppercase letter corresponds to phrases starting at the next position in the regular expression.

As an example from the area of programming languages, the set of all strings denoting a floating point number can be described by an extended right-regular grammar G with $N = \{S, A, B, C, D, E, F\}$, $\Sigma = \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9, +, -, ., e\}$, where S is the start symbol, and P consists of the following rules:

$S \rightarrow +A$	$A \rightarrow 0A$	$B \rightarrow 0C$	$C \rightarrow 0C$	$D \rightarrow +E$	$E \rightarrow 0F$	$F \rightarrow 0F$
$S \rightarrow -A$	$A \rightarrow 1A$	$B \rightarrow 1C$	$C \rightarrow 1C$	$D \rightarrow -E$	$E \rightarrow 1F$	$F \rightarrow 1F$
$S \rightarrow A$	$A \rightarrow 2A$	$B \rightarrow 2C$	$C \rightarrow 2C$	$D \rightarrow E$	$E \rightarrow 2F$	$F \rightarrow 2F$
	$A \rightarrow 3A$	$B \rightarrow 3C$	$C \rightarrow 3C$		$E \rightarrow 3F$	$F \rightarrow 3F$

$A \rightarrow 4A$	$B \rightarrow 4C$	$C \rightarrow 4C$	$E \rightarrow 4F$	$F \rightarrow 4F$
$A \rightarrow 5A$	$B \rightarrow 5C$	$C \rightarrow 5C$	$E \rightarrow 5F$	$F \rightarrow 5F$
$A \rightarrow 6A$	$B \rightarrow 6C$	$C \rightarrow 6C$	$E \rightarrow 6F$	$F \rightarrow 6F$
$A \rightarrow 7A$	$B \rightarrow 7C$	$C \rightarrow 7C$	$E \rightarrow 7F$	$F \rightarrow 7F$
$A \rightarrow 8A$	$B \rightarrow 8C$	$C \rightarrow 8C$	$E \rightarrow 8F$	$F \rightarrow 8F$
$A \rightarrow 9A$	$B \rightarrow 9C$	$C \rightarrow 9C$	$E \rightarrow 9F$	$F \rightarrow 9F$
$A \rightarrow .B$		$C \rightarrow eD$		$F \rightarrow \varepsilon$
$A \rightarrow B$		$C \rightarrow \varepsilon$		

Expressive power

There is a direct one-to-one correspondence between the rules of a (strictly) right-regular grammar and those of a nondeterministic finite automaton, such that the grammar generates exactly the language the automaton accepts.^[3] Hence, the right-regular grammars generate exactly all regular languages. The left-regular grammars describe the reverses of all such languages, that is, exactly the regular languages as well.

Every strict right-regular grammar is extended right-regular, while every extended right-regular grammar can be made strict by inserting new non-terminals, such that the result generates the same language; hence, extended right-regular grammars generate the regular languages as well. Analogously, so do the extended left-regular grammars.

If empty productions are disallowed, only all regular languages that do not include the empty string can be generated.^[4]

While regular grammars can only describe regular languages, the converse is not true: regular languages can also be described by non-regular grammars.

Mixing left-regular and right-regular rules

If mixing of left-regular and right-regular rules is allowed, we still have a linear grammar, but not necessarily a regular one. What is more, such a grammar need not generate a regular language: all linear grammars can be easily brought into this form, and hence, such grammars can generate exactly all linear languages, including non-regular ones.

For instance, the grammar G with $N = \{S, A\}$, $\Sigma = \{a, b\}$, P with start symbol S and rules

$$\begin{aligned} S &\rightarrow aA \\ A &\rightarrow Sb \\ S &\rightarrow \varepsilon \end{aligned}$$

generates $\{a^i b^i : i \geq 0\}$, the paradigmatic non-regular linear language.

See also

- Regular expression, a compact notation for regular grammars
- Regular tree grammar, a generalization from strings to trees
- Prefix grammar
- Chomsky hierarchy
- Perrin, Dominique (1990), "Finite Automata", in Leeuwen, Jan van (ed.), *Formal Models and Semantics*, Handbook of Theoretical Computer Science, vol. B, Elsevier, pp. 1–58
- Pin, Jean-Éric (Oct 2012). *Mathematical Foundations of Automata Theory* (<http://www.liafa.jussieu.fr/~jep/PDF/MPRI/MPRI.pdf>) (PDF)., chapter III

References

1. John E. Hopcroft and Jeffrey D. Ullman (1979). *Introduction to Automata Theory, Languages, and Computation* (<https://archive.org/details/introductiontoau00hopc>). Reading/MA: Addison-Wesley. ISBN 0-201-02988-X. Here: p.217 (left, right-regular grammars as subclasses of context-free grammars), p.79 (context-free grammars)
2. Hopcroft and Ullman 1979 (p.229, exercise 9.2) call it a normal form for right-linear grammars.
3. Hopcroft and Ullman 1979, p.218-219, Theorem 9.1 and 9.2
4. Hopcroft and Ullman 1979, p.229, Exercise 9.2

Retrieved from "https://en.wikipedia.org/w/index.php?title=Regular_grammar&oldid=1223384762"