

Definition: A context-free grammar $(V, A, \langle s \rangle, P)$ is called a regular grammar if every production rule in P is of one of the three forms:

- (i) $\langle A \rangle \rightarrow b \langle B \rangle$
- (ii) $\langle A \rangle \rightarrow b$
- (iii) $\langle A \rangle \rightarrow \varepsilon$

where $\langle A \rangle$ and $\langle B \rangle$ are nonterminals, b is a terminal, and ε is the empty word. A regular grammar is said to be in normal form if all its production rules are of types (i) and (iii).

Remark: In the literature, you often see this definition labelled left-regular grammar as opposed to right-regular grammar, where the production rules of types 1 have the form $\langle A \rangle \rightarrow \langle B \rangle b$, (**i.e.** the terminal is on the right of the nonterminal). This distinction is not really important as long as we stick to one type throughout since both left-regular grammars and right-regular grammars generate regular languages.

Lemma: Any language generated by a regular grammar may be generated by a regular grammar in normal form.

Proof: Let $\langle A \rangle \rightarrow b$ be a rule of type (ii). Replace it by two rules: $\langle A \rangle \rightarrow b \langle F \rangle$ and $\langle F \rangle \rightarrow \varepsilon$, where $\langle F \rangle$ is a new nonterminal. Add $\langle F \rangle$ to the set V . We do the same for every rule of type (ii) obtaining a bigger set V , but now our production rules are only of type (i) and (iii) and we are generating the same language.

qed

Example: Recall the regular language $L = \{0^m 1^n \mid m, n \in \mathbb{N}, m \geq 0, n \geq 0\}$. We can generate it from the regular grammar in normal form given by production rules:

- 1. $\langle s \rangle \rightarrow 0 \langle A \rangle$
- 2. $\langle A \rangle \rightarrow 0 \langle A \rangle$
- 3. $\langle A \rangle \rightarrow \varepsilon$
- 4. $\langle s \rangle \rightarrow \varepsilon$

5. $\langle A \rangle \rightarrow 1 \langle B \rangle$
6. $\langle B \rangle \rightarrow 1 \langle B \rangle$
7. $\langle s \rangle \rightarrow 1 \langle B \rangle$
8. $\langle B \rangle \rightarrow \varepsilon$

Rules (1), (2), (5), (6), (7) are of type (i), where rules (3), (4) and (8) are of types (iii).

(1) and (3) gives 0. (1), (2) applied $m - 1$ times and (3) gives 0^m for $m \geq 2$.

(7) and (8) give 1. (7), (6) applied $n - 1$ times and (8) give 1^n for $n \geq 2$. (1), (5) and (8) give 01. (1), (5), (6) applied $n - 1$ times and (8) gives 01^n for $n \geq 2$.

(1), (2) applied $m - 1$ times, (5) and (8) gives $0^m 1$ for $m \geq 2$.

(1), (2) applied $m - 1$ times, (5), (6) applied $n - 1$ times, and (8) gives $0^m 1^n$ for $m \geq 2, n \geq 2$.

Rule (4) gives the empty word $\varepsilon = 0^0 1^0$.

Q: Why does a regular grammar yield a regular language, **i.e.** one recognised by a finite state acceptor?

A: Not obvious from the definition, but we can construct the finite state acceptor from the regular grammar as follows: our regular grammar is given by $(V, A, \langle s \rangle, P)$. Want a finite state acceptor (S, A, i, t, F) . Immediately, we see the alphabet A is the same and $i = \langle s \rangle$. This gives us the idea of associating to every nonterminal symbol in $V \setminus A$ a state. $\langle s \rangle \in V \setminus A$, so that's good. Next we ask:

Q: Is it sufficient for $S = V \setminus A$?

A: No! Our set F of finishing/accepting states should be nonempty. So we add an element $\{f\}$ to $V \setminus A$, where our acceptor will end up when a word in our language. Thus, $S = (V \setminus A) \cup \{f\}$ and $F = \{f\}$. $F \subseteq S$ as needed.

Q: How do we define t ?

A: Use the production rules in P ! For every rule of type (i), which is of the form $\langle A \rangle \rightarrow b \langle B \rangle$ set $t(\langle A \rangle, b) = \langle B \rangle$. This works out well because our nonterminals $\langle A \rangle$ and $\langle B \rangle$ are states of the acceptor and the terminal $b \in A$ so t takes an element of $S \times A$ to an element of S as needed. Now look at production rules of type (ii), $\langle A \rangle \rightarrow b$ and of types (iii), $\langle A \rangle \rightarrow \varepsilon$. Those are applied when we finish constructing a word w in our language L , **i.e.** at the very last step, so our acceptor should end up in the finishing state f whenever a production rule of type (ii) or (iii) is applied. Write a production rule of type (ii) or (iii) as $\langle A \rangle \rightarrow w$, then we can set $t(\langle A \rangle, w) = f$. We have finished constructing t as well. Technically, $t : S \times (A \cup \{\varepsilon\}) \rightarrow S$ instead of $t : S \times A \rightarrow S$, but we can easily fix the transition function t by combining the last two transitions for each accepted word.

Remark: The same general principles as we used above allow us to go from a finite state acceptor to a regular grammar. This gives us the following theorem:

Theorem: A language L is regular $\Leftrightarrow L$ is recognised by a finite state acceptor $\Leftrightarrow L$ is generated by a regular grammar.

Regular expressions

11

Task Understand another equivalent way of characterizing regular languages due to Kleene in the 1950's.

Def. Let A be an alphabet.

1. \emptyset , ϵ , and all elements of A are regular expressions;
2. If w and w' are regular expressions, then $w \circ w'$, $w \cup w'$, and w^* are regular expressions.

Remark This definition is an inductive one.

NB It is important not to confuse the regular expressions \emptyset and ϵ . The expression ϵ represents the language containing a single string, namely ϵ the empty string, whereas \emptyset represents the language that doesn't contain any strings. Recall that a language L is any subset of $A^* = \bigcup_{n=0}^{\infty} A^n = A^0 \cup A^1 \cup A^2 \cup \dots$

$\{ \epsilon \}$ set of words of length 0
 A set of words of length 1
 A^2 set of words of length 2

Precedence order of operations if parentheses aren't present:
 First $*$, then \circ (concatenation), then \cup (union).

Examples

(1) $A = \{0, 1\}$ $1^*0 = \{ w \in A^* \mid w = 1^m 0 \text{ for } m \in \mathbb{N}, m \geq 0 \}$
 $= \{ 0, 10, 110, 1110, \dots \}$

1^*0 we can omit the concatenation symbol

(2) $A = \{a, b\}$ a^*b^* $(a^*b^*)^*$

$$(3) A = \{0, 1\} \quad (A \circ A)^* = \{w \in A^* \mid w \text{ is a word of even length}\}$$

Recall $L^* = \bigcup_{n=0}^{\infty} L^n$ where $L^0 = \{\epsilon\}$

$$L^1 = L \quad \text{and inductively}$$

$$L^n = L \circ L^{n-1}$$

$$\text{Here } L = A \circ A = \{00, 01, 10, 11\}$$

$$(3)' (A^* \circ A^*)^* = A^*$$

$$(4) A = \{0, 1\}$$

$$(0 \cup \epsilon) \circ (1 \cup \epsilon) = \{\epsilon, 0, 1, 01\}$$

$$(5) \epsilon^* = \{\epsilon\}$$

(6) $\phi^* = \{\epsilon\}$ The star operation concatenates any number of words from the language. If the language is empty, then the star operation can only put together 0 words, which yields only the empty word.

Use of regular expressions in programming

→ design of compilers for programming languages

Elemental objects in a programming language, which are called tokens

(for example variables names and constants) can be described w/

regular expressions → get the syntax of a programming language

this way. There exists an algorithm for recognizing regular expressions

that has been implemented \Rightarrow an automatic system generates the lexical

analyser that checks the input in a compiler.

→ eliminate redundancy in programming

The same regular expression can be generated in more than one

way (obvious from the definition of a regular expression) \Rightarrow

there exists an equivalence relation on regular expressions and

algorithms that check when two regular expressions are equivalent.