

Kleene's algorithm

In theoretical computer science, in particular in formal language theory, **Kleene's algorithm** transforms a given nondeterministic finite automaton (NFA) into a regular expression. Together with other conversion algorithms, it establishes the equivalence of several description formats for regular languages. Alternative presentations of the same method include the "elimination method" attributed to Brzozowski and McCluskey, the algorithm of McNaughton and Yamada,^[1] and the use of Arden's lemma.

Algorithm description

According to Gross and Yellen (2004),^[2] the algorithm can be traced back to Kleene (1956).^[3] A presentation of the algorithm in the case of deterministic finite automata (DFAs) is given in Hopcroft and Ullman (1979).^[4] The presentation of the algorithm for NFAs below follows Gross and Yellen (2004).^[2]

Given a nondeterministic finite automaton $M = (Q, \Sigma, \delta, q_0, F)$, with $Q = \{ q_0, \dots, q_n \}$ its set of states, the algorithm computes

the sets R_{ij}^k of all strings that take M from state q_i to q_j without going through any state numbered higher than k .

Here, "going through a state" means entering *and* leaving it, so both i and j may be higher than k , but no intermediate state may. Each set R_{ij}^k is represented by a regular expression; the algorithm computes them step by step for $k = -1, 0, \dots, n$. Since there is no state numbered higher than n , the regular expression R_{0j}^n represents the set of all strings that take M from its start state q_0 to q_j . If $F = \{ q_1, \dots, q_f \}$ is the set of accept states, the regular expression $R_{01}^n \mid \dots \mid R_{0f}^n$ represents the language accepted by M .

The initial regular expressions, for $k = -1$, are computed as follows for $i \neq j$:

$$R_{ij}^{-1} = a_1 \mid \dots \mid a_m \quad \text{where } q_j \in \delta(q_i, a_1), \dots, q_j \in \delta(q_i, a_m)$$

and as follows for $i = j$:

$$R_{ii}^{-1} = a_1 \mid \dots \mid a_m \mid \varepsilon \quad \text{where } q_i \in \delta(q_i, a_1), \dots, q_i \in \delta(q_i, a_m)$$

In other words, R_{ij}^{-1} mentions all letters that label a transition from i to j , and we also include ε in the case where $i = j$.

After that, in each step the expressions R_{ij}^k are computed from the previous ones by

$$R_{ij}^k = R_{ik}^{k-1} (R_{kk}^{k-1})^* R_{kj}^{k-1} \mid R_{ij}^{k-1}$$

Another way to understand the operation of the algorithm is as an "elimination method", where the states from 0 to n are successively removed: when state k is removed, the regular expression R_{ij}^{k-1} , which describes the words that label a path from state $i > k$ to state $j > k$, is rewritten into R_{ij}^k so as to take into account the

possibility of going via the "eliminated" state k .

By induction on k , it can be shown that the length^[5] of each expression R_{ij}^k is at most $\frac{1}{3}(4^{k+1}(6s+7) - 4)$ symbols, where s denotes the number of characters in Σ . Therefore, the length of the regular expression representing the language accepted by M is at most $\frac{1}{3}(4^{n+1}(6s+7)f - f - 3)$ symbols, where f denotes the number of final states. This exponential blowup is inevitable, because there exist families of DFAs for which any equivalent regular expression must be of exponential size.^[6]

In practice, the size of the regular expression obtained by running the algorithm can be very different depending on the order in which the states are considered by the procedure, i.e., the order in which they are numbered from 0 to n .

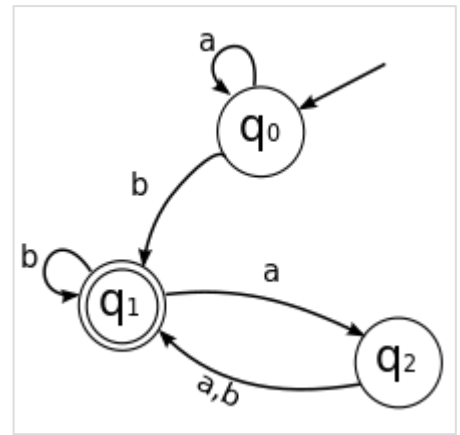
Example

The automaton shown in the picture can be described as $M = (Q, \Sigma, \delta, q_0, F)$ with

- the set of states $Q = \{q_0, q_1, q_2\}$,
- the input alphabet $\Sigma = \{a, b\}$,
- the transition function δ with $\delta(q_0, a) = q_0$, $\delta(q_0, b) = q_1$, $\delta(q_1, a) = q_2$, $\delta(q_1, b) = q_1$, $\delta(q_2, a) = q_1$, and $\delta(q_2, b) = q_1$,
- the start state q_0 , and
- set of accept states $F = \{q_1\}$.

Kleene's algorithm computes the initial regular expressions as

$$\begin{aligned} R_{00}^{-1} &= a \mid \varepsilon \\ R_{01}^{-1} &= b \\ R_{02}^{-1} &= \emptyset \\ R_{10}^{-1} &= \emptyset \\ R_{11}^{-1} &= b \mid \varepsilon \\ R_{12}^{-1} &= a \\ R_{20}^{-1} &= \emptyset \\ R_{21}^{-1} &= a \mid b \\ R_{22}^{-1} &= \varepsilon \end{aligned}$$



Example DFA given to Kleene's algorithm

After that, the R_{ij}^k are computed from the R_{ij}^{k-1} step by step for $k = 0, 1, 2$. Kleene algebra equalities are used to simplify the regular expressions as much as possible.

Step 0

$$R_{00}^0 = R_{00}^{-1} (R_{00}^{-1})^* R_{00}^{-1} \mid R_{00}^{-1} = (a \mid \varepsilon) (a \mid \varepsilon)^* (a \mid \varepsilon) \mid a \mid \varepsilon = a^*$$

$$\begin{aligned}
R_{01}^0 &= R_{00}^{-1} (R_{00}^{-1})^* R_{01}^{-1} \mid R_{01}^{-1} = (a \mid \varepsilon) (a \mid \varepsilon)^* b \mid b = a^* b \\
R_{02}^0 &= R_{00}^{-1} (R_{00}^{-1})^* R_{02}^{-1} \mid R_{02}^{-1} = (a \mid \varepsilon) (a \mid \varepsilon)^* \emptyset \mid \emptyset = \emptyset \\
R_{10}^0 &= R_{10}^{-1} (R_{00}^{-1})^* R_{00}^{-1} \mid R_{10}^{-1} = \emptyset \mid (a \mid \varepsilon)^* (a \mid \varepsilon) \mid \emptyset = \emptyset \\
R_{11}^0 &= R_{10}^{-1} (R_{00}^{-1})^* R_{01}^{-1} \mid R_{11}^{-1} = \emptyset \mid (a \mid \varepsilon)^* b \mid b \mid \varepsilon = b \mid \varepsilon \\
R_{12}^0 &= R_{10}^{-1} (R_{00}^{-1})^* R_{02}^{-1} \mid R_{12}^{-1} = \emptyset \mid (a \mid \varepsilon)^* \emptyset \mid a = a \\
R_{20}^0 &= R_{20}^{-1} (R_{00}^{-1})^* R_{00}^{-1} \mid R_{20}^{-1} = \emptyset \mid (a \mid \varepsilon)^* (a \mid \varepsilon) \mid \emptyset = \emptyset \\
R_{21}^0 &= R_{20}^{-1} (R_{00}^{-1})^* R_{01}^{-1} \mid R_{21}^{-1} = \emptyset \mid (a \mid \varepsilon)^* b \mid a \mid b = a \mid b \\
R_{22}^0 &= R_{20}^{-1} (R_{00}^{-1})^* R_{02}^{-1} \mid R_{22}^{-1} = \emptyset \mid (a \mid \varepsilon)^* \emptyset \mid \varepsilon = \varepsilon
\end{aligned}$$

Step 1

$$\begin{aligned}
R_{00}^1 &= R_{01}^0 (R_{11}^0)^* R_{10}^0 \mid R_{00}^0 = a^* b \mid (b \mid \varepsilon)^* \emptyset \mid a^* = a^* \\
R_{01}^1 &= R_{01}^0 (R_{11}^0)^* R_{11}^0 \mid R_{01}^0 = a^* b \mid (b \mid \varepsilon)^* (b \mid \varepsilon) \mid a^* b = a^* b^* b \\
R_{02}^1 &= R_{01}^0 (R_{11}^0)^* R_{12}^0 \mid R_{02}^0 = a^* b \mid (b \mid \varepsilon)^* a \mid \emptyset = a^* b^* ba \\
R_{10}^1 &= R_{11}^0 (R_{11}^0)^* R_{10}^0 \mid R_{10}^0 = (b \mid \varepsilon) (b \mid \varepsilon)^* \emptyset \mid \emptyset = \emptyset \\
R_{11}^1 &= R_{11}^0 (R_{11}^0)^* R_{11}^0 \mid R_{11}^0 = (b \mid \varepsilon) (b \mid \varepsilon)^* (b \mid \varepsilon) \mid b \mid \varepsilon = b^* \\
R_{12}^1 &= R_{11}^0 (R_{11}^0)^* R_{12}^0 \mid R_{12}^0 = (b \mid \varepsilon) (b \mid \varepsilon)^* a \mid a = b^* a \\
R_{20}^1 &= R_{21}^0 (R_{11}^0)^* R_{10}^0 \mid R_{20}^0 = (a \mid b) (b \mid \varepsilon)^* \emptyset \mid \emptyset = \emptyset \\
R_{21}^1 &= R_{21}^0 (R_{11}^0)^* R_{11}^0 \mid R_{21}^0 = (a \mid b) (b \mid \varepsilon)^* (b \mid \varepsilon) \mid a \mid b = (a \mid b) b^* \\
R_{22}^1 &= R_{21}^0 (R_{11}^0)^* R_{12}^0 \mid R_{22}^0 = (a \mid b) (b \mid \varepsilon)^* a \mid \varepsilon = (a \mid b) b^* a \mid \varepsilon
\end{aligned}$$

Step 2

$$\begin{aligned}
R_{00}^2 &= R_{02}^1 (R_{22}^1)^* R_{20}^1 \mid R_{00}^1 = a^* b^* ba \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid a^* = a^* \\
R_{01}^2 &= R_{02}^1 (R_{22}^1)^* R_{21}^1 \mid R_{01}^1 = a^* b^* ba \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid (a \mid b) b^* \mid a^* b^* b = a^* b (a (a \mid b) \mid b)^* \\
R_{02}^2 &= R_{02}^1 (R_{22}^1)^* R_{22}^1 \mid R_{02}^1 = a^* b^* ba \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid ((a \mid b) b^* a \mid \varepsilon) \mid a^* b^* ba = a^* b^* b (a (a \mid b) b^*)^* a \\
R_{10}^2 &= R_{12}^1 (R_{22}^1)^* R_{20}^1 \mid R_{10}^1 = b^* a \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid \emptyset = \emptyset \\
R_{11}^2 &= R_{12}^1 (R_{22}^1)^* R_{21}^1 \mid R_{11}^1 = b^* a \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid (a \mid b) b^* \mid b^* = (a (a \mid b) \mid b)^* \\
R_{12}^2 &= R_{12}^1 (R_{22}^1)^* R_{22}^1 \mid R_{12}^1 = b^* a \mid ((a \mid b) b^* a \mid \varepsilon)^* \mid ((a \mid b) b^* a \mid \varepsilon) \mid b^* a = (a (a \mid b) \mid b)^* a
\end{aligned}$$

$$\begin{aligned}
R_{20}^2 &= R_{22}^1 (R_{22}^1)^* R_{20}^1 \mid = ((a|b)b^*a \mid ((a|b)b^*a \mid \emptyset \mid \emptyset = \emptyset \\
R_{21}^2 &= R_{22}^1 (R_{22}^1)^* R_{21}^1 \mid = ((a|b)b^*a \mid ((a|b)b^*a \mid (a|b)b^* \mid (a \mid b) b^* = (a \mid b) (a (a \mid b) \mid b)^* \\
R_{22}^2 &= R_{22}^1 (R_{22}^1)^* R_{22}^1 \mid = ((a|b)b^*a \mid ((a|b)b^*a \mid ((a|b)b^*a \mid (a \mid b) b^* a \mid \varepsilon = ((a \mid b) b^* a)^*
\end{aligned}$$

Since q_0 is the start state and q_1 is the only accept state, the regular expression R_{01}^2 denotes the set of all strings accepted by the automaton.

See also

- Floyd–Warshall algorithm — an algorithm on weighted graphs that can be implemented by Kleene's algorithm using a particular Kleene algebra
- Star height problem — what is the minimum stars' nesting depth of all regular expressions corresponding to a given DFA?
- Generalized star height problem — if a complement operator is allowed additionally in regular expressions, can the stars' nesting depth of Kleene's algorithm's output be limited to a fixed bound?
- Thompson's construction algorithm — transforms a regular expression to a finite automaton

References

1. McNaughton, R.; Yamada, H. (March 1960). "Regular Expressions and State Graphs for Automata". *IRE Transactions on Electronic Computers*. **EC-9** (1): 39–47. doi:10.1109/TEC.1960.5221603 (https://doi.org/10.1109%2FTEC.1960.5221603). ISSN 0367-9950 (https://www.worldcat.org/issn/0367-9950).
2. Jonathan L. Gross and Jay Yellen, ed. (2004). *Handbook of Graph Theory*. Discrete Mathematics and it Applications. CRC Press. ISBN 1-58488-090-2. Here: sect.2.1, remark R13 on p.65
3. Kleene, Stephen C. (1956). "Representation of Events in Nerve Nets and Finite Automata" (http://www.dlsi.ua.es/~mlf/nnafmc/papers/kleene56representation.pdf) (PDF). *Automata Studies, Annals of Math. Studies*. **34**. Princeton Univ. Press. Here: sect.9, p.37-40
4. John E. Hopcroft, Jeffrey D. Ullman (1979). *Introduction to Automata Theory, Languages, and Computation* (https://archive.org/details/introductiontoau00hopc). Addison-Wesley. ISBN 0-201-02988-X. Here: Section 3.2.1 pages 91-96
5. More precisely, the number of regular-expression symbols, " a_i ", " ε ", "|", " $*$ ", " $.$ "; not counting parentheses.
6. Gruber, Hermann; Holzer, Markus (2008). "Finite Automata, Digraph Connectivity, and Regular Expression Size". In Aceto, Luca; Damgård, Ivan; Goldberg, Leslie Ann; Halldórsson, Magnús M.; Ingólfssdóttir, Anna; Walukiewicz, Igor (eds.). *Automata, Languages and Programming*. Lecture Notes in Computer Science. Vol. 5126. Springer Berlin Heidelberg. pp. 39–50. doi:10.1007/978-3-540-70583-3_4 (https://doi.org/10.1007%2F978-3-540-70583-3_4). ISBN 9783540705833. S2CID 10975422 (https://api.semanticscholar.org/CorpusID:10975422).. Theorem 16.

■