

Compiler: Note #1

ant-hengxin

0130

Big Picture:

we are used to expressing lexical rules in regular expressions, but we are writing the lexical analyzer through a DFA. That means that there exists some way to convert from regular expressions to DFA.

$$\text{re} \rightarrow \text{DFA}$$

But this approach is overly complicated, and for simplicity's sake, let's do this thing in multiple steps:

$$\text{re} \xrightarrow[\text{Construction}]{\text{Thompson}} \text{NFA} \xrightarrow[\text{Construction}]{\text{Subset}} \text{DFA}$$

And we also want to know how to convert a DFA to a re.

Alphabet Σ : A set of finite symbols.

String s over alphabet Σ : A sequence of symbols from Σ . A special string ϵ is the empty string with the property $|\epsilon| = 0$.

String operation: concatenation: $x = \text{dog}, y = \text{house} \Rightarrow xy = \text{doghouse}$.

Language L over alphabet Σ : A countable set of strings over Σ .

Language operation:

Suppose L, M are languages, we can use set operations to construct new languages:

1. $L \cup M := \{ s \mid s \in L \text{ or } s \in M \}$
2. $LM := \{ s \in L \text{ and } s \in M \}$
3. $L^* := \bigcup_{i=0}^{\infty} L^i$ (Kleene closure)
4. $L^+ := \bigcup_{i=1}^{\infty} L^i$ (Positive closure)

The following logic is, we firstly talk about the definition of regular expression, then we talk about the definition of non-deterministic finite automaton (NFA), and we define what's the meaning of a regular expression is equivalent to a NFA.

Regular expression:

A regular expression over alphabet Σ is defined as follows:

1. ϵ is a regular expression.
2. $\forall a \in \Sigma, a$ is a regular expression.
3. If r is a regular expression, then (r) is also a regular expression.
4. If r, s are regular expressions, then $r|s, rs, r^*$ are also regular expressions.

Regular expressions are prioritized as follows: $() > * > \text{concat} > |$.

Regular expressions define a language we call the regular language. The rules are as follows:

1. $L(\epsilon) = \{ \epsilon \}$.

2. $L(a) = \{ a \}, \forall a \in \Sigma$.
3. $L((r)) = L(r)$.
4. $L(r|s) = L(r) \cup L(s), L(rs) = L(r)L(s), L(r^*) = (L(r))^*$.

Note: If r_1, r_2 is a regular expression, then:

1. $r_1?$ represents 0 or 1 r_1 .
2. r_1/r_2 represents Indicates that r_1 when followed by r_2 .

NFA: non-deterministic-finite automaton:

A NFA is a five-tuple represented as $\mathcal{A} = (\Sigma, S, s_0, \delta, F)$ with the following properties:

1. Alphabet $\Sigma (\epsilon \notin \Sigma)$.
2. S is a finite set of states.
3. $s_0 \in S$ is the only start state.
4. $\delta : S \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^S$
5. F is a set of accept states.

Note that the requirement **only** in property 3 is not necessary, because if there exists more than one start state, we can define a new start state and the original start states is constructed as a ϵ -closure of the new start state.

And the definition of NFA actually define a language $L(\mathcal{A})$, which is composed of all the strings that can be accepted by the NFA.

Two important questions of NFA:

1. Given a string s , is s in $L(\mathcal{A})$?
2. What is $L(\mathcal{A})$?

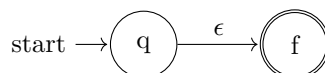
Equivalence of regular expression and NFA:

Now we can talk about what is the meaning of a regular expression is equivalent to a NFA. We call a regular expression r is equivalent to a NFA \mathcal{A} if and only if $L(r)$ is equivalent to $L(\mathcal{A})$ as two sets.

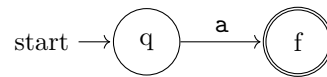
Let's start to talk about how to convert a **re** to a **NFA**, which let us introduce an algorithm called **Thompson's construction**.

So the idea of Thompson's construction is very simple, in order to know how to convert a regular expression to a NFA, what we only need to do is to know the conversion rules corresponding to the definition of regular expression.

1. if ϵ is a regular expression, then the corresponding NFA is:



2. if \mathbf{a} is a regular expression, then the corresponding NFA is:



3. The NFA of (\mathbf{s}) is the same as the NFA of \mathbf{s} .
4. The NFA of $\mathbf{r|s}$ (Given $\mathbf{NFA(s)}$ and $\mathbf{NFA(r)}$) is:

