

Derivatives of Regular Expressions

Robert Keller
 22 January 2007

There are many different notations for regular expressions. We work with the most basic version. Versions that are used in software tools, particularly various Unix tools, are more complex because they introduce many new symbols as short cuts, but no additional power.

Our Definition of Regular Expression

A regular expression over an alphabet Σ , and the language denoted by the expression is defined inductively as follows:

	Expression R	L(R), the language denoted by R	Interpretation
Basis	\emptyset	\emptyset	empty set
	ϵ	$\{\epsilon\}$	set of just the empty string
	σ where $\sigma \in \Sigma$	$\{\sigma\}$	set of one one-letter string
Induction rules	$(R + S)$	$L(R) \cup L(S)$	union
	(RS)	$L(R)L(S)$	concatenation
	$(R)^*$	$L(R)^*$	star

Inductive definitions permit proofs by structural induction, as in the construction of an NFA from any regular expression.

We often take short cuts in writing regular expressions, by omitting parentheses by using precedence rules. * binds more tightly than concatenation. Concatenation binds more tightly than union. For example, $abc^* + d$ abbreviates $((a (b (c^*))) + d)$.

Identities for Regular Expressions

These identities are derived from considering the interpretation of regular expressions. They are equalities not in that the expressions are the same, but in that the languages they denote are the same, just like many other algebraic identities. Here R, S, and T are arbitrary regular expressions.

1. $R + \emptyset = \emptyset + R = R$
2. $R\emptyset = \emptyset R = \emptyset$
3. $R\epsilon = \epsilon R = R$
4. $(R + S)T = RT + ST$
5. $R(S + T) = RS + RT$

6. $(R + S)^* = (R^*S^*)^*$
7. $(R^*)^* = R^*$
8. $R(ST) = (RS)T$
9. $R(SR)^* = (RS)^*R$
10. $R^n + R^* = R^*$ for any n , including $n = 0$ ($R^0 = \epsilon$), $n = 1$ ($R^1 = R$).
11. $R + S = S + R$
12. $R + (S + T) = (R + S) + T$
13. $(RS)^* = (R + S)^*$, provided that $\epsilon \in L(R)$.

Note that identities involving only $+$ (union) are known from set theory.

Derivatives of Regular Expressions

The idea of derivative of a regular expression can be used for constructing a DFA from an RE by hand, which is often more convenient than creating an NFA then performing the subset construction to get a DFA. The derivative of a regular expression R with respect to an alphabet symbol $\sigma \in \Sigma$, written R/σ , is defined by structural induction as follows:

	Expression	Derivative	Rule number
Basis	\emptyset/σ	\emptyset	1
	ϵ/σ	\emptyset	2
	σ/σ	ϵ	3
	σ/σ' where $\sigma' \neq \sigma$	\emptyset	4
Induction rules	$(RS)/\sigma$	$(R/\sigma)S$ if $\epsilon \notin L(R)$	5
		$(R/\sigma)S + S/\sigma$ if $\epsilon \in L(R)$	6
	$(R + S)/\sigma$	$(R/\sigma) + (S/\sigma)$	7
	$(R)^*/\sigma$	$(R/\sigma)R^*$	8

Derivatives provide a formal way to manipulate languages via regular expressions. The idea is that:

$$L(R/\sigma) = \{x \in \Sigma^* \mid \sigma x \in L(R)\}$$

In other words, $L(R/\sigma)$ is the set of all strings in $L(R)$ that begin with σ , but with the initial σ lopped off. (Strings that don't begin with σ do not contribute to $L(R/\sigma)$.)

Examples:

Suppose that $\Sigma = \{0, 1\}$. We have the following derivatives:

Directly from rules:

$$\begin{aligned} \emptyset/0 &= \epsilon/0 = \emptyset/1 = \epsilon/1 = 1/0 = 0/1 = \emptyset \\ 0/0 &= 1/1 = \epsilon \end{aligned}$$

Applying rules recursively, and using identities to simplify:

$$(00)/0 \text{ (using rule 5)} = (0/0) 0 = \epsilon 0 = 0$$

$$(01)/0 \text{ (using rule 5)} = (0/0) 1 = \epsilon 1 = 1$$

$$(00)/1 \text{ (using rule 5)} = (0/0) 0 = \epsilon 0 = 0$$

$$(01)/1 \text{ (using rule 5)} = (0/0) 1 = \epsilon 1 = 1$$

$$(0 + 1)/0 \text{ (using rule 7)} = 0/0 + 1/0 = \epsilon + \emptyset = \epsilon$$

$$(0 + 1)/1 \text{ (using rule 7)} = 0/1 + 1/1 = \emptyset + \epsilon = \epsilon$$

$$(0^*1)/0 = ((0^*/0)1 + 1/0) \text{ (using rule 6)} = ((0/0)0^*)1 + \emptyset = (\epsilon 0^*)1 = 0^*1$$

$$(0^*1)/1 = ((0^*/1)1 + 1/1) \text{ (using rule 6)} = ((0/1)0^*)1 + \epsilon = (\emptyset 0^*)1 + \epsilon = \emptyset 1 + \epsilon = \epsilon$$

Application

Derivatives of regular expressions can be used to construct a DFA from a regular expression without creating an NFA, as follows:

1. Simplify R first, if possible, as that will save work later.
2. Set New = {R}, Old = \emptyset .
3. While New $\neq \emptyset$
 - a. Let R be an arbitrary element of New. Move R from New to Old.
 - b. Let S be the derivatives of R with respect to each alphabet symbol using the above rules, simplifying where possible. Keep track of the expression from which each derivative is derived.
 - c. Add to New any simplified expressions not equivalent to an expression in New \cup Old.
4. Create a DFA with the regular expressions in Old as states.
5. The initial state is the original expression R.
6. There is a connection from R to R' labeled with σ exactly when R' is equivalent to R/ σ .
7. The accepting states are those corresponding to expressions that have ϵ in their language.

Example

$$R = (01)^*1 + 1$$

We can simplify this right away to $(01)^*1$, since $1 \in L((01)^*1)$.

$$\text{New} = \{(01)^*1\}$$

Move $(01)^*1$ from New to Old.

$$\begin{aligned} \text{Extract derivatives of } (01)^*1: \\ ((01)^*1)/0 &= ((01)^*/0)1 = (1(01)^*)1 \end{aligned}$$

$$((01)^*1)/1 = ((01)^*/1)1 + (1/1) = ((0/1)(01)^*)1 + (1/1) = \emptyset + \varepsilon = \varepsilon$$

$$\text{New} = \{(1(01)^*)1, \varepsilon\}$$

Extract derivatives of $(1(01)^*)1$:

$$\begin{aligned} (1(01)^*)1/0 &= \emptyset \\ (1(01)^*)1/1 &= (01)^*1, \text{ which is in Old.} \end{aligned}$$

Extract derivatives of ε :

$$\begin{aligned} \varepsilon/0 &= \emptyset, \text{ which is in Old.} \\ \varepsilon/1 &= \emptyset, \text{ which is in Old.} \end{aligned}$$

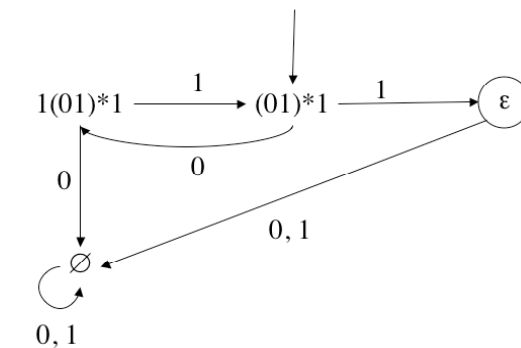
$$\text{New} = \{\emptyset\}$$

Extract derivatives of \emptyset :

$$\begin{aligned} \emptyset/0 &= \emptyset, \text{ which is in Old.} \\ \emptyset/1 &= \emptyset, \text{ which is in Old.} \end{aligned}$$

New is empty. Old = $\{(01)^*1, 1(01)^*1, \varepsilon, \emptyset\}$.

Drawn below is the DFA thus derived.



We can check visually that a regular expression representing its language is $(01)^*1$, as desired.