

*lec 5:Regular Expression*

## 1.Regular Expression

Regular Expression have an important role in computer science applications. in application involving text, users may want to search for strings that satisfy certain pattern. regular expression provided a powerful method for describing such pattern such as AWK and GREP in UNIX, Regular Expression is a set of symbols, Thus if alphabet= {a, b}, then aab, a, baba, bbbbb, and baaaaa would all be strings of symbols of alphabet.

**Regular expressions** can be used to define languages. A regular expression is like a "pattern"; strings that match the pattern are in the language, strings that do not match the pattern are not in the language. The construction of regular expressions is defined recursively, starting with primitive regular expressions, which can be composed using typical operators to form more complex regular expressions.

In addition we include an empty string denoted by (  $\epsilon$  ) which has no symbols in it.

## 2.Operation on Language

Regular operations on languages (sets of strings). Suppose  $L$  and  $K$  are languages.

- **Union:**  $L \cup K = \{x \mid x \in L \text{ or } x \in K\}$ .
- **Concatenation:**  $L \circ K = LK = \{xy \mid x \in L \text{ and } y \in K\}$ .
- **Star (Kleene star):**

$$L^* = \{w_1w_2 \dots w_n \mid w_1, \dots, w_n \in L \text{ and } n \geq 0\}.$$

We (hopefully) all understand what union does. The other two have some subtleties. Let

$$L = \{\text{under, over}\}, \quad \text{and} \quad K = \{\text{ground, water, work}\}.$$

Then

$$LK = \{\text{underground, underwater, underwork, overground, overwater, overwork}\}.$$

Similarly,

$$K^* = \left\{ \begin{array}{l} \epsilon, \text{ground, water, work, groundground,} \\ \text{groundwater, groundwork, workground,} \\ \text{waterworkwork, \dots} \end{array} \right\}.$$

### 3. Definition of Regular Expressions :

Let us fix an alphabet  $\Sigma$ . Here are the basic regular expression.

regex	conditions	set represented
$a$	$a \in \Sigma$	$\{a\}$
$\epsilon$		$\{\epsilon\}$
$\emptyset$		$\{\}$

Thus,  $\emptyset$  represents the empty language. But  $\epsilon$  represents that language which has the empty word as its only word in the language..

In particular, for a regular expression hexpi, we will use the notation  $L(\text{hexpi})$  to denote the language associated with this regular expression. Thus,

$$L(\epsilon) = \{ \epsilon \} \quad \text{and} \quad L(\emptyset) = \{ \} ;$$

which are two different languages.

Suppose that  $L(R)$  is the language represented by the regular expression  $R$ .

Here are recursive rules that make complex regular expressions out of simpler ones.

regex	conditions	set represented
$R \cup S$ or $R + S$	$R, S$ regexes	$L(R) \cup L(S)$
$R \circ S$ or $RS$	$R, S$ regexes	$L(R)L(S)$
$R^*$	$R$ a regex	$L(R)^*$

And some handy shorthand notation:

regex	conditions	set represented
$R^+$	$R$ a regex	$L(R)L(R)^*$
$\Sigma$		$\Sigma$

Some specific boundary case examples:

1.  $R\epsilon = R = \epsilon R$ .
2.  $R\emptyset = \emptyset = \emptyset R$ .

This is a bit confusing, so let us see why this is true, recall that

$$R\emptyset = \{xy \mid x \in R \text{ and } y \in \emptyset\}.$$

But the empty set ( $\emptyset$ ) does not contain any element, and as such, no concatenated string can be created. Namely, its the empty language.

3.  $R \cup \emptyset = R$  (just like with any set).
4.  $R \cup \epsilon = \epsilon \cup R$ .  
This expression can not always be simplified, since  $\epsilon$  might not be in the language  $L(R)$ .
5.  $\emptyset^* = \{\epsilon\}$ , since the empty word is always contain in the language generated by the star operator.
6.  $\epsilon^* = \{\epsilon\}$ .

### Examples of Kleene star:

$(1^*)$  is the set of strings  $\{\epsilon, 1, 11, 111, 1111, 11111, \text{etc.}\}$

$(1100)^*$  is the set of strings  $\{\epsilon, 1100, 11001100, 110011001100, \text{etc.}\}$

$(00+11)^*$  is the set of strings  $\{\epsilon, 00, 11, 0000, 0011, 1100, 1111, 000000, 000011, 001100, \text{etc.}\}$

$(0+1)^*$  is all possible strings of zeros and ones, often written as  $\Sigma^*$  where  $\Sigma = \{0, 1\}$

$(0+1)^* (00+11)$  is all strings of zeros and ones that end with either 00 or 11.

$(w)^+$  is a shorthand for  $(w)(w)^*$   $w$  is any string or expression and the superscript plus, +

### Concatenation:

Notation to the concatenation:  $.$  (The dot.):

if  $L1 = \{x, xxx\}$  and  $L2 = \{xx\}$  So  $(L1.L2)$  means  $L1$  concatenated  $L2$  and it is  $L1 \circ L2 = \{xxx, xxxxx\}$

Example 1:

$L1 = \{a, b\}$ .

$L2 = \{c, d\}$ .

$L1.L2 = \{ac, ad, bc, bd\}$

Note: ab differ from ba.

Example 2:

$\Sigma = \{x\}$ .

$L1 = \{\text{set of all odd words over } \Sigma \text{ with odd length}\}$ .

$L1 = \{\text{set of all even words over } \Sigma \text{ with odd length}\}$ .

$L1 = \{x, xxx, xxxxx, xxxxxxx, \dots\}$ .

$L2 = \{\_, xx, xxx, xxxxx, \dots\}$ .

$L1.L2 = \{x, xxx, xxxxx, xxxxxxx, \dots\}$ .

Example 3:

$L1 = \{x, xxx\}$ .

$L2 = \{xx\}$ .

$L1.L2 = \{xxx, xxxxx\}$ .

Some rules on concatenation:

$$\lambda.x = x$$

$L1.L2 = \{\text{set of elements}\}$

Example 4:

$A = \{a,b\}$  // the alphabet is composed of a and b

$A^* = \{1, a,b,aa,ab,ba,bb,aaa,aab, \dots\}$

The symbol \* is called the Kleene star.

$\emptyset$ (empty set)

$\epsilon$  (empty string)

Given regular expressions  $x$  and  $y$ ,  $x + y$  is a regular expression representing the set of all strings in either  $x$  or  $y$  (set union)

$$x = \{a b\} \quad y = \{c d\} \quad x + y = \{a b c d\}$$

Example 5:

Let  $A = \{0,1\}$ ,  $W1 = 00110011$ ,  $W2 = 00000$

$$W1W2 = 0011001100000$$

$$W2W1 = 0000000110011$$

$$W1\lambda = W1 = 00110011$$

$$\lambda W2 = W2 = 00000$$

$$x = \{a,b\} \quad y = \{c,d\} \quad xy = \{ac, ad, bc, bd\}$$

$$\text{Note: } (a + b)^* = (a^*b^*)^*$$

### More Examples of regular expressions

Describe the language = what is the output (words, strings) of the following RE

Regular expression	output(set of strings)
$\lambda$	$\{\lambda\}$
$\lambda^*$	$\{\lambda\}$
$a$	$\{a\}$
$aa$	$\{aa\}$
$a^*$	$\{\epsilon, a, aa, aaa, \dots\}$
$aa^*$	$\{a, aa, aaa, \dots\}$
$a^+$	$\{a, aa, aaa, \dots\}$
$ba^+$	$\{ba, baa, baaa, \dots\}$
$(ba)^+$	$\{ba, baba, bababa, \dots\}$
$(a b)$	$\{a, b\}$

$a b^*$	$\{ a, l, b, bb, bbb, \dots \}$
$(a b)^*$	$\{\lambda, a, b, aa, ab, ba, bb, \dots \}$
$aa(ba)^*bb$	$\{ aabb, aababb, aabababb, \dots \}$
$(a + a)$	$\{ a \}$
$(a + b)$	$\{ a, b \}$
$(a + b)^2$	$(a + b)(a + b) == \{ aa, ab, ba, bb \}$
$(a + b + c)$	$\{ a, b, c \}$
$(a + b)^*$	$\{ \epsilon, a, b, aa, bb, ab, ba, aaa, bbb, aab, bba, \dots \}$
$(abc)$	$\{ abc \}$
$(\lambda + a) bc$	$\{ bc, abc \}$
$ab^*$	$\{ a, ab, abb, abbb, \dots \}$
$(ab)^*$	$\{\lambda, ab, abab, ababab, \dots \}$
$a + b^*$	$\{ a, \lambda, b, bb, bbb, \dots \}$
$a (a + b)^*$	$\{ a, aa, ab, aaa, abb, aba, abaa, \dots \}$ (a
$+ b)^* a (a + b)^*$	$\{ a, aaa, aab, baa, bab, \dots \}$
$(a + \lambda)^*$	$(a)^* = \{\lambda, a, aa, aaa, \dots \}$
$x^* (a + b) + (a + b)$	$x^* (a + b)$
$x^* y + y$	$x^* y$
$(x + 1)x^*$	$x^* (x + 1) = x^*$
$(x + 1)(x + 1)^* (x + 1)$	$x^*$

note: symbol  $\lambda$  same mean epsilon symbol ( $\epsilon$ )