

Regular Expression

(Lecture 2)

A Regular Expression can be recursively defined as follows:

1. ϵ is a Regular Expression indicating the language containing an empty string. ($L(\epsilon) = \{\epsilon\}$ one string, with no symbols).
2. ϕ is a Regular Expression denoting an empty language. ($L(\phi) = \{\}$ no strings).
3. x is a Regular Expression where $L = \{x\}$ one string, with one symbol x .
4. If X is a Regular Expression denoting the language $L(X)$ and Y is a Regular Expression denoting the language $L(Y)$, then:
 - a) $X + Y$ is a Regular Expression corresponding to the language $L(X) \cup L(Y)$ where $L(X+Y) = L(X) \cup L(Y)$.
 - b) $X \cdot Y$ is a Regular Expression corresponding to the language $L(X) \cdot L(Y)$ where $L(X \cdot Y) = L(X) \cdot L(Y)$
 - c) R^* is a Regular Expression corresponding to the language $L(R^*)$ where $L(R^*) = (L(R))^*$
5. If we apply any of the rules several times from 1 to 5, they are Regular Expressions.

- $L(R_1 \cup R_2) = L(R_1) \cup L(R_2)$
- $L(R_1 \circ R_2) = L(R_1) \circ L(R_2)$
- $L(R_1^*) = (L(R_1))^*$

Example

Expression $((0 \cup 1) \epsilon)^* \cup 0$ denotes language $\{0, 1\}^* \cup \{0\} = \{0, 1\}^*$, all strings.

Example

$(0 \cup 1)^* 111 (0 \cup 1)^*$ denotes $\{0, 1\}^* \{111\} \{0, 1\}^*$, all strings with substring 111.

Example

$L =$ strings over $\{0, 1\}$ with odd number of 1s. $0^* 1 0^* (0^* 1 0^* 1 0^*)^*$

Example

$L =$ strings with substring 01 or 10.

$$[(0 \cup 1)^* 01 (0 \cup 1)^*] \cup [(0 \cup 1)^* 10 (0 \cup 1)^*]$$

Abbreviate (writing Σ for $(0 \cup 1)$):

$$\Sigma^* 01 \Sigma^* \cup \Sigma^* 10 \Sigma^*$$

Example

L = strings with substring 01 or 10.

$$(0 \cup 1)^* 01 (0 \cup 1)^* \cup (0 \cup 1)^* 10 (0 \cup 1)^*$$

Abbreviate:

$$\Sigma^* 01 \Sigma^* \cup \Sigma^* 10 \Sigma^*$$

Example

L = strings with neither substring 01 or 10.

- Can't write complement.
- But can write: $0^* \cup 1^*$.

Example

L = strings with no more than two consecutive 0s or two consecutive 1s

- Would be easy if we could write complement.

$$(\epsilon \cup 1 \cup 11) \cdot ((0 \cup 00) \cdot (1 \cup 11))^* (\epsilon \cup 0 \cup 00)$$

- Alternate one or two of each.

- Regular expressions commonly used to specify syntax.
- For (portions of) programming languages –Editors – Command languages like UNIX shell

Example

Decimal numbers $DD^* \cdot D^* \cup D^* \cdot DD^*$,

Where D is the alphabet $\{0, \dots, 9\}$ Need a digit either before or after the decimal point.

Some RE Examples

Regular Expression	Regular Set
$(0+10^*)$	$L = \{0, 1, 10, 100, 1000, 10000, \dots\}$
(0^*10^*)	$L = \{1, 01, 10, 010, 0010, \dots\}$
$(0+\epsilon).(1+\epsilon)$	$L = \{\epsilon, 0, 1, 01\}$
$(a+b)^*$	Set of strings of a's and b's of any length including the null string. So $L = \{\epsilon, a, b, aa, ab, bb, ba, aaa, \dots\}$
$(a+b)^*abb$	Set of strings of a's and b's ending with the string abb. So $L = \{abb, aabb, babb, aaabb, ababb, \dots\}$
$(11)^*$	Set consisting of even number of 1's including empty string, So $L = \{\epsilon, 11, 1111, 111111, \dots\}$
$(aa)^*(bb)^*b$	Set of strings consisting of even number of a's followed by odd number of b's, So $L = \{b, aab, aabbb, aabbbbb, aaaab, aaaabbb, \dots\}$
$(aa + ab + ba + bb)^*$	String of a's and b's of even length can be obtained by concatenating any combination of the strings aa, ab, ba and bb including null, so $L = \{\epsilon, aa, ab, ba, bb, aaab, aaba, \dots\}$

- The languages that are associated with these regular expressions are called regular languages.

Example

Consider the language L Where $L = \{\Lambda x xx xxx \dots\}$ by using star notation, we may write $L = \text{language}(x^*)$.

Since x^* is any string of x's (including Λ).

Example

If we have the alphabet $\Sigma = \{a,b\}$ And $L = \{a ab abb abbb abbbb \dots\}$ Then $L = \text{language}(ab^*)$

Example

$(ab)^* = \Lambda$ or ab or abab or ababab or abababab or

Example

$L_1 = \text{language } (xx^*)$

The language L_1 can be defined by any of the expressions:

xx^* or x^+ or xx^*x^* or x^*xx^* or x^+x^* or x^*x^+ or $x^*x^*x^*xx^*$... Remember x^* can always be Λ .

Example

Language $(ab^*a) = \{aa\ aba\ abba\ abbba\ abbbbba\ \dots\}$

Example

Language $(a^*b^*) = \{\Lambda\ a\ b\ aa\ ab\ bb\ aaa\ aab\ abb\ bbb\ \dots\}$ ba and aba are not in this language so $a^*b^* \neq (ab)^*$

Example

The following expressions both define the language $L_2 = \{x^{\text{odd}}\}$: $x\ (xx)^*$ or $(xx)^*x$ But the expression x^*xx^* does not since it includes the word $(xx)\ x(x)$.

Example

Consider the language T defined over the alphabet $\Sigma = \{a,b,c\}$ $T = \{a\ c\ ab\ cb\ abb\ cbb\ abbb\ cbbb\ abbbb\ cbbbb\ \dots\}$

Then $T = \text{language } ((a+c)\ b^*)$ $T = \text{language } (\text{either } a \text{ or } c \text{ then some } b\text{'s})$

Example

Consider a finite language L that contains all the strings of a 's and b 's of length exactly three.

$L = \{aaa\ aab\ aba\ abb\ baa\ bab\ bba\ bbb\}$ $L = \text{language } ((a+b)\ (a+b)\ (a+b))$
 $L = \text{language } ((a+b)\ 3)$

Note: from the alphabet $\Sigma = \{a,b\}$, if we want to refer to the set of all possible strings of a 's or b 's of any length (including Λ) we could write $(a+b)^*$

Example

We can describe all words that begins with a and end with b with the expression $a\ (a+b)^*\ b$ which mean a (arbitrary string) b

Example

If we have the expression $(a+b)^* a (a+b)^*$ then the word abbaab can be considered to be of this form in three ways: $(\Lambda) a (bbaab)$ or $(abb) a (ab)$ or $(abba) a (b)$

Example

$(a+b)^* a (a+b)^* a (a+b)^* =$ (some beginning) (the first important a) (some middle) (the second important a) (some end)

Another expressions that denote all the words with at least two a's are: $b^* ab^* a (a+b)^*$, $(a+b)^* ab^* ab^*$, $b^* a (a+b)^* ab^*$

Then we could write:

$$\begin{aligned} &= \text{language } ((a+b)^* a (a+b)^* a (a+b)^*) \\ &= \text{language } (b^* ab^* a (a+b)^*) \\ &= \text{language } ((a+b)^* ab^* ab^*) \\ &= \text{language } (b^* a (a+b)^* ab^*) \\ &= \text{all words with at least two a's.} \end{aligned}$$

Note: two regular expressions are equivalent if they describe the same language.

Example

If we want all the words with exactly two a's, we could use the expression: $b^* ab^* ab^*$ which describe such words as aab, baba, bbbabbabbbb,...

Example

The language of all words that have at least one a and at least one b is

$$(a+b)^* a (a+b)^* b (a+b)^* + (a+b)^* b (a+b)^* a (a+b)^*$$

Note: $(a+b)^* b (a+b)^* a (a+b)^* \neq bb^* aa^*$ since the left includes the word aba, which the expression on the right side does not.

Note:

$$\begin{aligned} (a+b)^* &= (a+b)^* + (a+b)^* (a+b)^* = (a+b)^* (a+b)^* \\ (a+b)^* &= a (a+b)^* + b (a+b)^* + \Lambda \\ (a+b)^* &= (a+b)^* ab (a+b)^* + b^* a^* \end{aligned}$$

Note: usually when we employ the star operation we are defining an infinite language. We can represent a finite language by using the plus alone.

Example

$$L = \{abba\ baaa\ bbbb\}$$

$$L = \text{language}(abba + baaa + bbbb)$$

Example

$$L = \{\Lambda\ a\ aa\ bbb\}$$

$$L = \text{language}(\Lambda + a + aa + bbb)$$

Example

$$L = \{\Lambda\ a\ b\ ab\ bb\ abb\ bbb\ abbb\ bbbb\ \dots\}$$

We can define L by using the expression $b^* + ab^*$

Definition

The set of regular expressions is defined by the following rules:

Rule1: every letter of Σ can be made into a regular expression, Λ is a regular expression.

Rule2: if R_1 and R_2 are regular expressions, then so are: $(R_1) R_1R_2 R_1+R_2 R_1^*$.

Rule3: nothing else is a regular expression. Remember that $R_1^+ = R_1R_1^*$

Definition

If **S** and **T** are sets of strings of letters (whether they are finite or infinite sets), we define the product set of strings of letters to be: **ST** = {all combination of a string from **S** concatenated with a string from **T**}

Example

$$\text{If } S = \{a\ aa\ aaa\} \quad T = \{bb\ bbb\}$$

$$\text{Then } ST = \{abb\ abbb\ aabb\ aabbb\ aaabb\ aaabbb\} \quad (a+aa+aaa) (bb+bbb) \\ = abb+abbb+ aabb+aabbb+aaabb+aaabbb)$$

Example

If $P = \{a\} \cup \{bb\} \cup \{bab\}$ $Q = \{\Lambda\} \cup \{bbbb\}$ Then $PQ = \{a\} \cup \{bb\} \cup \{bab\} \cup \{abbbb\} \cup \{bbbbbb\} \cup \{babbbb\}$
 $(a+bb+bab)(\Lambda+bbbb) = a+bb+bab+ab^4+b^6+bab^5$

Example

If $M = \{\Lambda\} \cup \{x\} \cup \{xx\}$ $N = \{\Lambda\} \cup \{y\} \cup \{yy\} \cup \{yyy\} \cup \{yyyy\} \dots$

Then $MN = \{\Lambda\} \cup \{y\} \cup \{yy\} \cup \{yyy\} \cup \{yyyy\} \dots$
 $x\ xy\ xyy\ xyyy\ xyyyy\ \dots$
 $xx\ xxy\ xxyy\ xxyyy\ xxyyyy\ \dots$

Using regular expression we could write: $(\Lambda+x+xx)(y^*) = y^* + xy^* + xxy^*$

Definition

The following rules define the language associated with any regular expression.

Rule1: the language associated with the regular expression that is just a single letter is that one-letter word alone and the language associated with Λ is just $\{\Lambda\}$, a one-word language.

Rule2: if R_1 is regular expression associated with the language L_1 and R_2 is regular expression associated with the language L_2 then:

- i. The regular expression $(R_1)(R_2)$ is associated with the language L_1 times L_2 .

$$\text{Language } (R_1R_2) = L_1L_2$$

- ii. The regular expression R_1+R_2 is associated with the language formed by the union of the sets L_1 and L_2 .

$$\text{Language } (R_1+R_2) = L_1+L_2$$

- iii. The language associated with the regular expression $(R_1)^*$ is L_1^* , the kleene closure of the set L_1 as a set of words.

$$\text{Language } (R_1)^* = L_1^*$$

Example

$$L = \{\text{baa abba bababa}\}$$

The regular expression for this language is (baa+abba+bababa)

Example

$$L = \{\Lambda \text{ x xx xxx xxxxx xxxxxx}\}$$

The regular expression for this language is $(\Lambda+x+xx+xxx+xxxx+xxxxx)$

$$= (\Lambda+x)^5$$

Example

$$L = \text{language } ((a+b)^*(aa+bb)(a+b)^*)$$

$$= (\text{arbitrary}) (\text{double letter}) (\text{arbitrary})$$

$\{\Lambda \text{ a b ab ba aba bab abab baba } \dots\}$ these words are not included in L but they included by the regular expression: $(\Lambda+b)(ab)^*(\Lambda+a)$

Example

$$E = (a+b)^* a (a+b)^* (a+\Lambda) (a+b)^* a (a+b)^*$$

$$E = (a+b)^* a (a+b)^* a (a+b)^* a (a+b)^* + (a+b)^* a (a+b)^* \Lambda (a+b)^* a (a+b)^*$$

$$\text{We have } (a+b)^* \Lambda (a+b)^* = (a+b)^*$$

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

The language associated with E is not different from the language associated with: $(a+b)^* a (a+b)^* a (a+b)^*$

Note: $(a+b^*)^* = (a+b)^* (a^*)^* = a^* (aa+ab^*)^* \neq (aa+ab)^* (a^*b^*)^* = (a+b)^*$

Example

$$E = [aa+bb+(ab+ba)(aa+bb)^*(ab+ba)]^*$$

$$\text{Even-even} = \{\Lambda \text{ aa bb aabb abab abba baab baba bbaa aaaabb aaabab} \dots\}$$