

REGULAR EXPRESSIONS AND FINITE-STATE AUTOMATA

1)

- Let the alphabet $\Sigma = \{a, b\}$.
- a. Define a language L_1 over Σ to be the set of all strings that begin with the character **a** and have length at most three characters. Find L_1 .
- b. A **palindrome** is a string that looks the same if the order of its characters is reversed. For instance, aba and baab are palindromes. Define a language L_2 over Σ to be the set of all palindromes obtained using the characters of Σ . Write ten elements of L_2 .

Solution

- a. $L_1 = \{a, aa, ab, aaa, aab, aba, abb\}$
- b. L_2 contains the following ten strings (among infinitely many others):
 $\epsilon, a, b, aa, bb, aaa, bab, abba, babaabab, abaabbbbbaaba$

2)

- Let Σ be an alphabet. For each nonnegative integer n , let
- Σ^n = the set of all strings over Σ that have length n ,
- Σ^+ = the set of all strings over Σ that have length at least 1, and
- Σ^* = the set of all strings over Σ .
- Let $\Sigma = \{a, b\}$.
- a. Find Σ^0 , Σ^1 , Σ^2 , and Σ^3 .
- b. Let $A = \Sigma^0 \cup \Sigma^1$ and $B = \Sigma^2 \cup \Sigma^3$. Use words to describe A , B , and $A \cup B$.
- c. Describe a systematic way of writing the elements of Σ^+ . What change needs to be made to obtain the elements of Σ^* ?

Solution

- a. $\Sigma^0 = \{\epsilon\}$, $\Sigma^1 = \{a, b\}$, $\Sigma^2 = \{aa, ab, ba, bb\}$, and $\Sigma^3 = \{aaa, aab, aba, abb, baa, bab, bba, bbb\}$

- b. A is the set of all strings over Σ of length at most 1.

B is the set of all strings over Σ of length 2 or 3.

$A \cup B$ is the set of all strings over Σ of length at most 3.

- c. Elements of Σ^+ can be written systematically by writing all the strings of length 1, then all the strings of length 2, and so forth.

Σ^+ : a, b, aa, ab, ba, bb, aaa, aab, aba, abb, baa, bab, bba, bbb, aaaa,...

Of course the process of writing the strings in Σ^+ would continue forever, because Σ^+ is an infinite set. The only change that needs to be made to obtain Σ^* is to place the null string at the beginning of the list. ■

Activate
Go to Setti

3)

- In 1 and 2 let $\Sigma = \{x, y\}$ be an alphabet.

1. a. Let L_1 be the language consisting of all strings over Σ that are palindromes and have length ≤ 4 . List the elements of L_1 between braces.

- b. Let L_2 be the language consisting of all strings over Σ that begin with an x and have length ≤ 3 . List the elements of L_2 .

$$L_1 = \{x, y, xx, yy, xxx, yyy, xyx, yxy, xxxx, yyyy, xyxx, yxyx\}$$

$$L_2 = \{x, xx, xy, xxx, xxy, xyx, xyy\}$$

4)

- b. List between braces the elements of Σ^4 , the set of strings of length 4 over Σ .

$$\Sigma^4 = \{xxxx, xxxy, xyxx, xxyy, xyxy, xyxx, xyxy, xyxy, yxxx, yxxy, yxyx, yxyy, yyyx, yyxy, yyyy\}$$

Polish Notation

(a + b)_ **infix notation**_ operator such as + sits between the two quantities.

(+ab)_ **prefix notation**_ binary operator precedes the quantities on which it acts

(ab+)_ **postfix notation**_ binary operator follows the quantities on which it acts.

A great advantage of these notations is that they **eliminate the need for parentheses** in writing arithmetic expressions.

Ex: For instance, in postfix (or reverse Polish) notation, the expression $8 + 6 /$ is evaluated from left to right (Same as $(8+4)/6$).

Infix: $(a + b) \cdot c$ Postfix: $ab + c \cdot$

- a. If the expression $ab \cdot cd \cdot +$ in postfix notation is converted to infix notation, what is the result?
- b. Let $\Sigma = \{4, 1, +, -\}$, and let L = the set of all strings over Σ obtained by writing either a 4 or a 1 first, then either a 4 or a 1, and finally either a + or a -. List all elements of L between braces, and evaluate the resulting expressions

$$L = \{44+, 44-, 41+, 41-, 14+, 14-, 11+, 11-\}$$

Evaluation: $4+4=8$, $4-4=0$, etc.

LANGUAGES CAN BE COMBINED TO FORM NEW LANGUAGES

Let Σ be an alphabet. Given any strings x and y over Σ .

The **concatenation** of L and L' , denoted LL' , is: $LL' = \{xy \mid x \in L \text{ and } y \in L'\}$

The **union** of L and L' , denoted $L \cup L'$, is: $L \cup L' = \{x \mid x \in L \text{ or } x \in L'\}$

The **Kleene closure** of L , denoted L^* , is:

$L^* = \{x \mid x \text{ is a concatenation of any finite number of strings in } L\}$

Note **that ϵ is in L^*** because it is **regarded as a concatenation of zero strings** in L .

1)

- Let L_1 be the set of all strings consisting of an even number of a's (namely, ϵ , aa , $aaaa$, $aaaaaa$, \dots), and let $L_2 = \{b, bb, bbb\}$.

Find L_1L_2 , $L_1 \cup L_2$, and $(L_1 \cup L_2)^*$. Note that the null string ϵ is in L_1 because 0 is an even number.

I) L_1L_2 = the set of all strings that consist of an even number of a's followed by b or by bb or by bbb .

II) $L_1 \cup L_2$ = the set that includes the strings b , bb , bbb and any strings consisting of an even number of a's.

III) $(L_1 \cup L_2)^*$ = the set of all strings of a's and b's in which every occurrence of a is in a block consisting of an even number of a's.

The Language Defined by a Regular Expression

The regular expression r^* is called the Kleene closure of r

- Given an alphabet Σ , the following **are regular expressions over Σ** :
 - I. **BASE**: \emptyset, ϵ , and each individual symbol in Σ are regular expressions over Σ .
 - II. **RECURSION** : If r and s are regular expressions over Σ , then the following are also regular expressions over Σ :
 - (i) (rs) (ii) $(r \mid s)$ (iii) (r^*)

If the alphabet Σ happens to include symbols—such as $()^*$ —special provisions have to be made to **avoid ambiguity**. An escape character, usually a **backslash**, is added before the potentially ambiguous symbol

For instance, a left parenthesis would be written as $\backslash($ and the backslash itself would be written as $\backslash\backslash$.

To eliminate parentheses, **an order of precedence** for the operations used to define regular expressions has been introduced.

The **highest is $*$** , **concatenation** is next, and $|$ is the lowest.

It is also customary to **eliminate the outer set of parentheses** in a regular expression, because doing so does not produce ambiguity.

$$(a((bc)^*)) = a(bc)^* \text{ and } (a|(bc)) = a|bc$$

- a. Add parentheses to make the order of precedence clear in the following expression: $ab^* | b^*a$.
- b. Use the convention about order of precedence to eliminate the parentheses in the following expression: $((a|((b^*)c))(a^*))$.

Solution

- a. $((a(b^*))|((b^*)a))$
- b. $(a|b^*c)a^*$

Language defined by regular expression

- For any finite alphabet Σ , the function L that associates a language to each regular expression over Σ is defined by (I) and (II) below. For each such regular expression r , $L(r)$ is called the **language defined by r** .
- I. **BASE**: $L(\emptyset) = \emptyset$, $L(\epsilon) = \{\epsilon\}$, $L(a) = \{a\}$ for every a in Σ .
- II. **RECURSION**: If $L(r)$ and $L(r')$ are the languages defined by the regular expressions r and r' over Σ , then
 - (i) $L(rr') = L(r)L(r')$ (ii) $L(r|r') = L(r) \cup L(r')$ (iii) $L(r^*) = (L(r))^*$

Example: Let $\Sigma = \{a, b\}$, and consider the language defined by the regular expression $(a|b)^*$. Use set notation to find this language, and describe it in words.

$L((a|b)^*) = L(a|b)^* = L(\{a\} \cup \{b\})^* = \{a, b\}^* = \Sigma^*$ = set of all strings that start by characters a or b .

Note that **concatenating strings** and **taking unions** of sets are both associative operations. Thus for any regular expressions r, s and t ,

$L((r\ s)t) = L(r(st))$, Since $L((r|s)|t) = L(r|(s|t))$

The expression a^* matches only strings consisting entirely of the character a , including the empty string (e.g., "", "a", "aa", "aaa", etc.)

- In each of (a) and (b), let $\Sigma = \{a, b\}$ and consider the language L over Σ defined by the given regular expression.
- a. The regular expression is $a^*b(a \mid b)^*$. Write five strings that belong to L .
- b. The regular expression is $a^* \mid (ab)^*$. Indicate which of the following strings belong to L :

a b aaaa abba ababab

Solution

- a. The strings b , ab , $abbb$, $abaaa$, and $ababba$ are five strings from the infinitely many in L .
- b. The following strings are the only ones listed that belong to L : a , $aaaa$, and $ababab$. The string b does not belong to L because it is neither a string of a 's nor a string of possibly repeated ab 's. The string $abba$ does not belong to L because any two b 's that might occur in a string of L are separated by an a .

• Let $\Sigma = \{0, 1\}$. Use words to describe the languages defined by the following regular expressions over Σ .

- a. $0^*1^* \mid 1^*0^*$
 - b. $0(0 \mid 1)^*$
- a. The strings in this language consist either of a string of 0's followed by a string of 1's or of a string of 1's followed by a string of 0's. However, in either case the strings could be empty, which means that ϵ is also in the language.
 - b. The strings in this language have to start with a 0. The 0 may be followed by any finite number (including zero) of 0's and 1's in any order. Thus the language is the set of all strings of 0's and 1's that start with a 0. ■

