# Theory of Automata Regular Expressions

Hafiz Tayyeb Javed

Week 2 Lecture 2

# Contents

- Regular Expressions
  - Language Defining Symbols
  - Alternation, Either/OR, Disjunction, Plus Sign
- Defining Languages by Another New Method
- Formal Definition of Regular Expressions
  - Product Set
  - Languages Associated with Regular Expressions
  - Finite Languages Are Regular
  - How Hard It Is to Understand a Regular Expression
- Introducing EVEN-EVEN

# Regular Expressions

- RE is the sequence of characters or symbols that represent a finite or infinite set of text strings.

- *Pattern-matching* is the process of checking whether a text string conforms to a set of characteristics defined by patterns such as regular expressions.

- A regular expression is a set of pattern matching rules encoded in a string according to certain syntax rules. Although the syntax is somewhat complex it is very powerful and allows much more useful pattern matching than say simple wildcards like ? and *.

# Regular Expression

- A regular expression (sometimes abbreviated to "regex") is a way for a computer user or programmer to express how a computer program should look for a specified pattern in text and then what the program is to do when each pattern match is found.

- For example, a regular expression could tell a program to search for all text lines that contain the word "Windows 95" and then to print out each line in which a match is found or substitute another text sequence (for example, just "Windows") where any match occurs.

- The best known tool for specifying and handling the incidence of regular expressions is grep, a utility found in Unix-based operating systems and also offered as a separate utility program for Windows and other operating systems.

# Language-Defining Symbols

- We now introduce the use of the Kleene star, applied not to a set, but directly to the letter x and written as a superscript: x*.
- This simple expression indicates some sequence of x's (may be none at all):

$$\mathbf{x^*} = \Lambda \text{ or } x \text{ or } x^2 \text{ or } x^3 \ldots$$
$$= x^n \text{ for some } n = 0, 1, 2, 3, \ldots$$

- Letter **x** is intentionally written in boldface type to distinguish it from an alphabet character.

- We can think of the star as an unknown power. That is, **x*** stands for a string of x's, but we do not specify how many, and it may be the null string .

# R.E. Continued...

- The notation x* can be used to define languages by writing, say $L_4$ = language (x*)

- Since x* is any string of x's, $L_4$ is then the language of all possible strings of x's of any length (including Λ).

- *We should not confuse x* (which is a **language-defining symbol**) with $L_4$ (which is the **name** we have given to a certain language).*

# R.E. Continued...

- Given the alphabet = {a, b}, suppose we wish to define the language L that contains all words of the form: one *a* followed by some number of *b*'s (maybe no *b*'s at all); that is

- L = {a, ab, abb, abbb, abbbb, ...}

- Using the language-defining symbol, we may write

<div align="center">

L = language (ab*)

</div>

- This equation obviously means that L is the language in which the words are the concatenation of an initial a with some or no b's.

- *From now on, for convenience, we will simply say **some** b**'s** to mean **some or no** b**'s**. When we want to mean **some positive number of** b**'s**, we will explicitly say so.*

# R.E. Continued…

- We can apply the Kleene star to the whole string ab if we want:

    (ab)* = Λ  or ab or abab or ababab…

- Observe that

    (ab)* ≠ a*b*

    – because the language defined by the expression on the left contains the word abab, whereas the language defined by the expression on the right does not.

# R.E. Continued…

- If we want to define the language L1 = {x, xx, xxx, …} using the language-defining symbol, we can write

  L1 = language(xx*)

  which means that each word of L1 must start with an x followed by some (or no) x's.

- Note that we can also define L1 using the notation + (as an exponent) introduced in Chapter 2:

  L1 = language($x^+$)

  – which means that each word of L1 is a string of some positive number of x's.

# Alternation, Either/OR, Disjunction, Plus Sign

- Let us introduce another use of the plus sign. By the expression

    x + y

  where x and y are strings of characters from an alphabet, we mean **either** x **or** y.


- Care should be taken so as not to confuse this notation with the notation + (as an exponent) or with sign for arithmetic addition.

# Example

- Consider the language T over the alphabet

    $\Sigma = \{a; b; c\}$:

- T = {a; c; ab; cb; abb; cbb; abbb; cbbb; abbbb; cbbbb; …}

- In other words, all the words in T begin with either an a or a c and then are followed by some number of b's.

- Using the above plus sign notation, we may write this as

    T = language((a+ c)b*)

# Example

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

  L = {aaa, aab, aba, abb, baa, bab, bba, bbb}

- Note that the first letter of each word in L is either an a or a b; so are the second letter and third letter of each word in L.

- Thus, we may write

    L = language((a+ b)(a + b)(a + b))

- or for short,

    L = language((a+ b)$^3$)

# Example

- In general, if we want to refer to the set of all possible strings of a's and b's of any length whatsoever, we could write

  language((a+ b)*)

- This is the set of **all possible strings** of letters from the alphabet Σ = {a, b}, **including the null string**.

- This is powerful notation. For instance, we can describe all the words that begin with first an *a*, followed by anything (i.e., as many choices as we want of either a or b) as

  a(a + b)*

# Formal Definition of Regular Expressions

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

- **Rule 1:** Every letter of the alphabet $\Sigma$ can be made into a regular expression by writing it in **boldface, $\Lambda$** itself is a regular expression.

- **Rule 2:** If $r_1$ and $r_2$ are regular expressions, then so are:
    - **(i)** $(r_1)$
    - **(ii)** $r_1 r_2$
    - **(iii)** $r_1 + r_2$
    - **(iv)** $r_1*$

- **Rule 3:** Nothing else is a regular expression.

- Note: If $r_1 = aa + b$ then when we write $r_1*$, we really mean $(r_1)*$, that is $r_1* = (r_1)* = (aa + b)*$

# Example

- Consider the language defined by the expression

    (a + b)*a(a + b)*

- At the beginning of any word in this language we have

  *(a + b)\**, which is any string of *a*'s and *b*'s, then comes an *a*, then another any string.

- For example, the word abbaab can be considered to come from this expression by 3 different choices:

    *(Λ)a(bbaab)  or (abb)a(ab)     or (abba)a(b)*

# Example contd.

- This language is the set of all words over the alphabet Σ = {a, b} that have at least one a.

- The only words left out are those that have only b's and the word Λ.

  These left out words are exactly the language defined by the expression b*.

- If we combine this language, we should provide a language of all strings over the alphabet Σ = {a, b}. That is,

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

# Example

- Write RE to define the language of all words that have at least two a's :

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

- Another expression that defines all the words with at least two a's is

$$b*ab*a(a + b)*$$

- Hence, we can write

$$(a + b)*a(a + b)*a(a + b)* = b*ab*a(a + b)*$$

where by the equal sign we mean that these two expressions are **equivalent** in the sense that they describe the same language.

# Example

- The language of all words that have at least one **a** and at least one **b** is somewhat trickier. If we write

    (a + b)*a(a + b)*b(a + b)*

  then we are requiring that an *a* must precede a *b* in the word. Such

  words as ba and bbaaaa are not included in this language.


- Since we know that either the *a* comes before the *b* or the *b* comes before the *a*, we can define the language by the expression

    (a + b)*a(a + b)*b(a + b)* + (a + b)*b(a + b)*a(a + b)*


- Note that the only words that are omitted by the first term

  (a + b)*a(a + b)*b(a + b)* are the words of the form some b's followed by some a's. They are defined by the expression bb*aa*

# Example

- We can add these specific exceptions. So, the language of all words over the alphabet Σ = {a, b} that contain at least one **_a_** and at least one **_b_** is defined by the expression:

    (a + b)a(a + b)b(a + b) + bb*aa*

- Thus, we have proved that

(a + b)*a(a + b)*b(a + b)* + (a + b)*b(a + b)*a(a + b)*

= (a + b)*a(a + b)*b(a + b)* + bb*aa*

# Example

- In the above example, the language of all words that contain both an a and a b is defined by the expression

  (a + b)*a(a + b)*b(a + b)* + bb*aa*

- The only words that do not contain both an a and a b are the words of all a's, all b's, or Λ.

- When these are included, we get everything. Hence, the expression

  (a + b)*a(a + b)*b(a + b)* + bb*aa* + a* + b*

  defines all possible strings of a's and b's, including  (accounted for in both a and b).

- Thus

(a + b)* = (a + b)*a(a + b)*b(a + b)* + bb*aa* + a* + b*

# Example

- The following equivalences show that we should not treat expressions as algebraic polynomials:

    (a + b)* = (a + b)* + (a + b)*

    (a + b)* = (a + b)* + a*

    (a + b)* = (a + b)*(a + b)*

    (a + b)* = a(a + b)* + b(a + b)* + Λ

    (a + b)* = (a + b)*ab(a + b)* + b*a*

- The last equivalence may need some explanation:
    – The first term in the right hand side, (a + b)*ab(a + b)*, describes all the words that contain the substring ab.
    – The second term, b*a* describes all the words that do not contain the substring ab (i.e., all a's, all b's, Λ, or some b's followed by some a's).

# Example

- Let V be the language of all strings of a's and b's in which either the strings are all b's, or else an a followed by some b's. Let V also contain the word Λ. Hence,

  V = {Λ, a, b, ab, bb, abb, bbb, abbb, bbbb, …}

- We can define V by the expression

  b* + ab*

  where Λ is included in b*.

- Alternatively, we could define V by

  (Λ + a)b*

which means that in front of the string of some b's, we have

either an a or nothing.

# Example contd.

- Hence,

  (Λ + a)b* = b* + ab*


- Since b* = Λb*, we have

  (Λ + a)b* = b* + ab*

  which appears to be distributive law at work.


- However, we must be extremely careful in applying distributive law. Sometimes, it is difficult to determine if the law is applicable.

# Product Set

- 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 combinations of a string from S

                       concatenated with a string from T in that order}

# Example

- If S = {a, aa, aaa} and T = {bb, bbb} then

ST = {abb, abbb, aabb, aabbb, aaabb, aaabbb}

- – Note that the words are not listed in lexicographic order.

- Using regular expression, we can write this example as

(a + aa + aaa)(bb + bbb)
    = abb + abbb + aabb + aabbb + aaabb + aaabbb

# Example

- If M = {Λ, x, xx} and N = {Λ, y, yy, yyy, yyyy, …} then

- MN ={Λ, y, yy, yyy, yyyy,…x, xy, xyy, xyyy, xyyyy, …xx, xxy, xxyy, xxyyy, xxyyyy, …}

- Using regular expression

$$(\Lambda + x + xx)(y^*) = y^* + xy^* + xxy^*$$