

# Lecture - 18

## Topic: Regular Expressions

Scribed by: Sanskar Shaurya (22B0985)

Checked and compiled by:

**Disclaimer.** Please note this document has not received the usual scrutiny that formal publications enjoy. This may be distributed outside this class only with the permission of the instructor.

## 1 Regular Expressions

We know that NFAs with  $\epsilon$  are equivalent to NFAs without  $\epsilon$  which are equivalent to DFAs. Now we will introduce another way of defining a language known as **Regular Expressions** which are much easier to write.

### 1.1 Syntax of Regular Expressions

Let us say our alphabet  $\Sigma$  is  $\{a, b\}$ .

Our syntax will include two special symbols  $a, b$ . These are our basic building blocks alongside the symbol  $\epsilon$ .

If  $e_1$  and  $e_2$  are two regular expressions, then the following will also be regular expression:

- $e_1 + e_2$
- $e_1 \cdot e_2$
- $e_1^*$
- $(e_1)$

### 1.2 Semantics of Regular Expressions

Semantics of a Regular Expression will be a language.

- $\llbracket a \rrbracket = \{b\}$     $\llbracket b \rrbracket = \{b\}$     $\llbracket \epsilon \rrbracket = \{\epsilon\}$
- $\llbracket e_1 + e_2 \rrbracket = \llbracket e_1 \rrbracket \cup \llbracket e_2 \rrbracket$
- $\llbracket e_1 \cdot e_2 \rrbracket = \llbracket e_1 \rrbracket \cdot \llbracket e_2 \rrbracket = \{u \cdot v \mid u \in \llbracket e_1 \rrbracket, v \in \llbracket e_2 \rrbracket\}$

**Example:**  $\llbracket (a \cdot b) + a \rrbracket = \{ab, a\}$

The order of precedence among these operators is  $*$  then  $\cdot$  then  $+$ .

So the above regular expression is equivalent to  $a \cdot b + a$

- $e_1^n = \underbrace{e_1 \cdot e_1 \cdots e_1}_{n \text{ times}}$

Suppose  $e_1 = a + b$

Then  $\llbracket e_1 \rrbracket = \{a, b\}$  and  $\llbracket e_1^3 \rrbracket = \{aaa, aab, aba, abb, baa, bab, bba, bbb\}$

$\llbracket e_1^0 \rrbracket$  is used to denote  $\{\epsilon\}$ , the set containing just the empty string.

$$- \llbracket e_1^* \rrbracket = \bigcup_{i \geq 0} \llbracket e_1^i \rrbracket$$

**Examples:**

- $\llbracket (a+b)^* \rrbracket = \{\epsilon, a, b, aa, ab, ba, bb, \dots\}$   
 $\llbracket (a+b)^* \rrbracket$  is equivalent to  $\Sigma^*$  as it represents all finite length strings formed from  $\Sigma = \{a, b\}$
- $\llbracket a^* + b^* \rrbracket = \{u \in \Sigma^* \mid u = a^n \text{ or } b^n, n \geq 0\}$
- $\llbracket (a^*) \cdot (b^*) \rrbracket = \{u \in \Sigma^* \mid u = a^n \cdot b^m, n, m \geq 0\}$

**Q)** Is  $\llbracket ((a^*) \cdot (b^*))^* \rrbracket = L((a+b)^*)^1$ ?

Easy to see that  $\llbracket ((a^*) \cdot (b^*))^* \rrbracket \subseteq L((a+b)^*)$  as the language defined by the second regular expression contains every finite length string formed by  $\{a, b\}$ .

**Claim:**  $L((a+b)^*) \subseteq \llbracket ((a^*) \cdot (b^*))^* \rrbracket$

**Proof:** Consider any arbitrary string of length n, and then consider  $((a^*) \cdot (b^*))^n$ . As we are concatenating n strings from  $((a^*) \cdot (b^*))$ , we can pick a single letter *a* or *b* from each of the n strings corresponding to the letter at that position. As we can do this for any string with a positive length n, this means that every finite length string is in the language defined by the regular expression. *a* is by definition in  $\llbracket ((a^*) \cdot (b^*))^* \rrbracket$  and hence our claim is correct.

This means that  $\llbracket ((a^*) \cdot (b^*))^* \rrbracket = L((a+b)^*)$ . Hence just looking at syntax, we cannot decide whether the language defined by them is different or not, but we can do this easily after converting the regular expression to a DFA.

---

<sup>1</sup> $L(e_1)$  denotes the Language defined by the regular expression  $e_1$