# Lecture 5

Ben Rosenberg

June 14, 2021

**Note**: In these notes, the symbol ∘ is used where · may have been used in lecture. They mean the same thing (concatenation).

## Regular Languages and Regular Expressions

### Regular operations

Definition: **Regular operations** are applied to sets (languages). There are three, which follow:

(Note: $\Sigma = \{$a,b,c$\}$ for all examples below.)

#### Union

Union: The union operation is the usual one as applied to sets.

Example: Consider $L_1 = \{$a, aa, b, cbca, aba, ab, c$\}$ and $L_2 = \{\lambda,$a, b, cccba, cc$\}$. Then, the union of $L_1$ and $L_2$ is defined as $L_1 \cup L_2 = \{x | x \in L_1 \lor x \in L_2\}$. We say that $b \in L_1 \cup L_2$, because $b$ is in either $L_1$ or $L_2$ (the fact that it is in both is not relevant).

Union is commutative: that is, $L_1 \cup L_2 = L_2 \cup L_1$ for all languages $L_1$ and $L_2$.

#### Concatenation

Concatenation: The concatenation $L_1 \circ L_2 = \{w | (\exists x \in L_1)(\exists y \in L_2)(w = xy)\}$. It is the set of all strings that can be obtained by "gluing together" one element of the first language with one element of the second language.

Example: ab is an element of $L_1 \circ L_2$, because $a \in L_1$ and $b \in L_2$. Furthermore, $ab \in L_1$ and $\lambda \in L_2$. Either of these "splits" suffices to show that ab $\in L_1 \circ L_2$ – the quantifier is existential, meaning that we are able to choose whichever split is most useful to us.

Note that concatenation is *not* commutative. Consider the string ca, which is in $L_1 \circ L_2$. This string is not in $L_2 \circ L_1$ as none of the three possible splits of ca can be written as the concatenation of a string from $L_2$ and a string from $L_1$.

Assuming we have an algorithm to answer $x \in L_1$, or $x \in L_2$, then we have an algorithm to answer whether $x \in L_1 L_2$ or $x \in L_1 \circ L_2$. We simply take $x$ and try all possible "splits". For each split, we take the left side and check whether $L_1$ contains it, and take the right side, and check whether $L_2$ contains it. If we find a split for which the left side is in $L_1$ and the right side is in $L_2$, then we have shown that $x \in L_1 \circ L_2$.

Aside, before proceeding to the last operation:

**Cardinalities**   The cardinality of $|L_1 \cup L_2| \leq |L_1| + |L_2|$. For the above example, we have $|L_1 \cup L_2| = 7 + 5 - 2 = 10$ because $2$ of the elements are duplicates.

The union of $L_1$ and $L_2$ can be empty if both $L_1$ and $L_2$ are empty.

The cardinality of $|L_1 \circ L_2| \leq |L_1| \cdot |L_2|$, as there may be duplicate strings that can be created from the concatenation. ab, for example, belongs to $L_1 \circ L_2$ in two different ways, because a is in $L_1$ and b is in $L_2$, and ab is in $L_1$ and $\lambda$ is in $L_2$.

If either $L_1$ or $L_2$ is empty, then their concatenation $L_1 \circ L_2$ is empty – you can never find a string that is in a given language if it is empty, meaning that there are no valid concatenated strings.

Onto the third operation.

**Kleene star**

Notation: For any set $L$: $L^0 = \{\lambda\}$. Recursively, $L^{n+1} = L^n \circ L$.

For example, $L^1 = L \circ L^0 = L \circ \{\lambda\} = L$. Similarly, $L^2 = L \circ L = LL$.

Definition: The **Kleene star** of a language $L$, denoted $L^*$, is equal to:

$$L^* = \bigcup_{x \geq 0} L^k = \boxed{L^0 \cup L^1 \cup L^2 \cup \cdots}$$

This is, in effect, the set of all strings that can be obtained by gluing together zero or more strings of the original language.

Another example: aaa $\in L_1^*$. The split a|aa $\in L_1^2$ works, as do a|a|a $\in L_1^3$ and aa|a $\in L_1^2$.

In fact, there are no strings (from $\Sigma$ defined above) that are not contained in $L_1^*$ because each of the elements of $\Sigma$ are contained in $L_1^*$.

We define $\Sigma^*$ as the set of all strings over $\Sigma$. Since $\Sigma \subseteq L_1$, $\Sigma^* \subseteq L_1^*$.

The Kleene star of any language – even $\emptyset$, the empty language – always contains at least one element: namely, $\lambda$. So, $(\ )^* \neq \emptyset$ in all cases.

There are two cases in which $L^*$ is finite. The first is that in which $L = \emptyset$, which gives $\emptyset^* = \lambda$. The second is $L = \{\lambda\}$, which also gives $L^* = \{\lambda\}$.

If $L^*$ is not finite, then its cardinality is equal to $\aleph_0$. (See Gödel injection.)

## Regular expressions

Definition: Given an alphabet $\Sigma = \{a,b,c\}$, the **class of regular languages** over $\Sigma$ contains exactly those languages which are obtained by finitely many applications of regular operations to the sets:

$$\emptyset, \{\lambda\}, \{a\}, \{b\}, \{c\}$$

Definition: A **regular set** is any set which can be obtained by starting from the above, and applying any of the regular operations.

Definition: A **regular expression** over an alphabet $\{\Sigma = \{a,b,c\}\}$ is a string over the alphabet:

$$\{a, \ b, \ c, \ (, \ ), \ \cup, \ \circ, \ *, \ \emptyset, \ \lambda\}$$

This is the *programmer's alphabet* (our "keyboard"), which represents a regular language. Note that these are all *symbols* and do not actually mean anything until they are used to write regular expressions.

Observe that the alphabet "keyboard" used to define regular expressions contains $\Sigma$. It also contains the seven symbols that are not contained in $\Sigma$ and are used to represent notation.

Examples of translation from language to regular expression:

| language | regular expression |
|---|---|
| $\emptyset$ | $\emptyset$ |
| $\{\lambda\}$ | $\lambda$ |
| $\{a\}$ | a |
| recursively, if: | |
| $L_1 \mapsto$ | $\mapsto e_1$ |
| $L_2 \mapsto$ | $\mapsto e_2$ |
| then, | |
| $L_1 \cup L_2$ | $(e_1) \cup (e_2)$ |
| $L_1 \circ L_2$ | $(e_1) \circ (e_2)$ |
| $L_1^*$ | $(e_1)^*$ |

Order of operations:

1. $*$
2. $\circ$
3. $\cup$

Examples ( $\iff$ is used to denote going between regex and languages):

$$a \cup bc^* = (a) \cup \left( (b) \circ \left( (c)^* \right) \right)$$

$$a + bc^2 = (a) \cup \left( (b) \circ \left( (c)^2 \right) \right)$$

$$\emptyset \circ \lambda \iff \emptyset$$

$$a \iff \{a\}$$

$$a \cup b \iff \{a, b\}$$

$$ab \iff \{ab\}$$

$$(a \cup b)(a \cup b) = (a \cup b) \circ (a \cup b) \iff \{aa, ab, ba, bb\}$$

$$a^* \iff \{\lambda, a, aa, aaa, \dots \}$$

The assignment in class is always the following:

> Construct a regular expression that defines the set of exactly those strings over $\Sigma$ that satisfy the following property:

$$\text{consists of a's and b's} \implies (a \cup b)^*$$

Aside: "$A$ contains $B$" means that $B$ is a subset of or element of $A$. "$A$ consists of $B$" means that all of $A$ are defined as being $B$.

**More examples**

all possible strings

$\implies$ $(a \cup b \cup c)^*$

strings that have length equal to 3 (should be 27 of them)

$\implies$ $(a \cup b \cup c) \circ (a \cup b \cup c) \circ (a \cup b \cup c)$

length $\leq 3$

$\implies$ $(a \cup b \cup c \cup \lambda) \circ (a \cup b \cup c \cup \lambda) \circ (a \cup b \cup c \cup \lambda)$

consist of even # of a's

$\implies$ $(aa)^*$

Note that $(aa)^*$ means an even number of a's, while $aa*$ is any number of a's greater than 1.

even length

$\implies$ $\Big((a \cup b \cup c) \circ (a \cup b \cup c)\Big)^*$

odd length

$\implies$ $\Big((a \cup b \cup c) \circ (a \cup b \cup c)\Big)^* \circ (a \cup b \cup c)$

contains odd number of a's (contains means it could also have b's and c's)

$\implies$ $\big(a(b \cup c)^* a \cup b \cup c)^*\big)a(b \cup c)^*$

contains even number of a's (contains means it could also have b's and c's)

$\implies$ $(b \cup c)^* a(b \cup c)^* a(b \cup c)^*)^*(b \cup c)^*$

contain substring `abcb`

$\implies$ $(a \cup b \cup c)^* \mathtt{abcb} (a \cup b \cup c)^*$

contain as a substring either `abb` or `bca`

$\implies$ $(a \cup b \cup c)^* (\mathtt{abb} \cup \mathtt{bca})(a \cup b \cup c)^*$

contain as a substring both `abb` and `bca`

$\implies$ $(a \cup b \cup c)^* \mathtt{abb}(a \cup b \cup c)^* \mathtt{bca}(a \cup b \cup c)^* \bigcup (a \cup b \cup c)^* \mathtt{bca}(a \cup b \cup c)^* \mathtt{abb}(a \cup b \cup c)^* \bigcup (a \cup b \cup c)^* \mathtt{abbca}(a \cup b \cup c)^* \bigcup (a \cup b \cup c)^* \mathtt{bcabb}(a \cup b \cup c)^*$

evidently, intersection is a challenge.

begin with `bca`

$\implies$ $\mathtt{bca}(a \cup b \cup c)^*$

do not begin with `bca`

$\implies$ ... ?

We want, here, to find a means of a "complement". This one is the complement of the previous answer. We could write out the 26 correct 3-letter strings that are not `bca`, and then add $(a \cup b \cup c)^*$ to the end, except that the length does not have to be 3 or greater.

We know that all strings of length 2 or less are good. Given the assumption that all strings of length 2 or less are good, then we can actually do it without enumerating all 27 of them.

We have:

$$(a \cup b \cup c \cup \lambda)^*(a \cup b \cup c \cup \lambda)^*$$

$$\bigcup$$

$$(a \cup c)(a \cup b \cup c)^*$$

$$\bigcup$$

$$b(a \cup b)(a \cup b \cup c)^*$$

$$\bigcup$$

$$bc(b \cup c)(a \cup b \cup c)^*$$

Theorem (to be proven later): If there is a regular expression for some language, then there is a regular expression for its complement.

Aside: Good book: "Regular Algebra and Finite Machines", by John Conway