# Regular Expressions and Languages

*Prof. Licia Sbattella*

*aa 2007-08*

*Translated and adapted by L. Breveglieri*

# REGULAR EXPRESSIONS AND LANGUAGES

REGULAR LANGUAGES are the simplest family
of formal languages. This family can be defined
in different ways;
- algebraically
- by means of generating grammars
- by means of recognizer machines

$$\Sigma = \{a_1, a_2, \dots a_i\}$$

$$\cdot \quad \cup \quad *$$

$$\{a_1\}, \{a_2\}, \dots \{a_i\} \ \varnothing$$

REGULAR EXPRESSION (regexp) is a string $r$ over the letters
of the terminal alphabet $\Sigma$ and containing a few metasymbols,
according to the following rules (where $s$ and $t$ represent
regular expressions):

metasymbols

$$\cdot \quad \cup \quad * \quad \varnothing \quad (\ )$$

1. $r = \varnothing$     3. $r = (s \cup t)$    5. $r = (s)*$

2. $r = a, \ a \in \Sigma$    4. $r = (s.t) \ o \ r = (st)$

$P$ay attention: U is often
denoted | (vertical bar)

PRECEDENCE OF OPERATORS: star ($*$), concatenation (.), union (U)

In the regexp it is permitted to use ε, because it holds: $\ \varepsilon = \varnothing^*$

THE MEANING OF A REGEXP *e* is a language $L_e$ over the alphabet Σ according to the following table:

| | regexp | language $L_r$ |
|---|---|---|
| 1. | $\varepsilon$ | $\{\varepsilon\}$ |
| 2. | $a \in \Sigma$ | $\{a\}$ |
| 3. | $s \bigcup t \text{ o } s \mid t$ | $L_s \bigcup L_t$ |
| 4. | $s.t \text{ o } st$ | $L_s . L_t$ |
| 5. | $s *$ | $L_s^{*}$ |

**A LANGUAGE IS REGULAR**

if it is the meaning
(= if it is generated)
by a regular expression

EXAMPLE 1: $L_e$ is the (unary) language of the multiples of three

$$e = (111) *$$
$$L_e = \{\varepsilon, 111, 111111, ....\} = \{1^n \mid n \bmod 3 = 0\}$$
$$e_1 = 11(1) * \qquad L_e \neq L_{e_1}$$
$$L_{e_1} = \{11, 111, 1111, 11111, ....\} = \{111^n \mid n \geq 0\}$$

EXAMPLE 2: Let $\Sigma = \{ +, -, d \}$, where $d$ denotes a decimal digit 0,1,…,9.

Define the regexp $e$ that generates the language of integers, possibly signed.

$$e = (+ \cup - \cup \varepsilon)dd*$$

$$L_e = \{+, -, \varepsilon\}\{d\}\{d\}^*$$

EXAMPLE 3: define a language over the binary alphabet { a, b }, such that the number of letters $a$ in each string is odd and that each string contains at least one letter $b$.

$$A_p b A_d \mid A_d b A_p$$
$$A_p = b^*(ab^*ab^*)^* \qquad A_d = b^*ab^*(ab^*ab^*)^*$$

FAMILY OF FINITE LANGUAGES (FIN):
is the collection of all languages of finite cardinality
(= that contain only finitely many strings).

Caution: REG and FIN are sets of sets of strings,
not sets of strings themselves.

EVERY FINITE LANGUAGE IS REGULAR, because it is the (finite)
union of finitely many strings, and each string is in turn the concatenation
of finitely many letters or in general of alphabetical symbols.

$$(x_1 \bigcup x_2 \bigcup ... \bigcup x_k) = (a_{1_1} a_{1_2} ... a_{1_n} \bigcup ... \bigcup a_{k_1} a_{k_2} ... a_{k_m})$$

But the family of regular languages also contains
many (actually infinitely many) languages of infinite
cardinality (= that contain infinitely many strings).

$$FIN \subset REG$$

1. consider e regexp that is fully parenthesized
2. number every alphebetical symbol occurring in the regexp
3. isolate the subexpressions and put them into evidence

$$e = (a \cup (bb))^* (c^+ \cup (a \cup (bb)))$$

$$e_N = (a_1 \cup (b_2 b_3))^* (c_4{}^+ \cup (a_5 \cup (b_6 b_7)))$$

$(a_1 \cup (b_2 b_3))^*$      $c_4{}^+ \cup (a_5 \cup (b_6 b_7))$

$a_1 \cup (b_2 b_3)$      $c_4{}^+$    $a_5 \cup (b_6 b_7)$

$a_1$    $b_2 b_3$      $c_4$    $a_5$    $b_6 b_7$

$b_2$   $b_3$          $b_6$   $b_7$

# DIFFERENT PARENTHESIZATIONS (and hence interpretations) OF A REGEXP

$$f = (a \cup bb)^* (c^+ \cup a \cup bb)$$

$$(a \cup bb)^* (c^+ \cup (a \cup bb))$$

$$(a \cup bb)^* ((c^+ \cup a) \cup bb)$$

$$
\begin{array}{llll}
(a_1 \cup b_2 b_3)^* & c_4^+ \cup a_5 \cup b_6 b_7 & & \\
a_1 \cup b_2 b_3 & c_4^+ & a_5 \cup b_6 b_7 & c_4^+ \cup a_5 \quad b_6 b_7 \\
a_1 \quad b_2 b_3 & c_4 & a_5 \quad b_6 b_7 & \\
\quad b_2 \quad b_3 & & b_6 \quad b_7 &
\end{array}
$$

CHOICE: the union and iteration operators that occur in a regexp correspond to some (sometimes infinitaly many) possible alternatives.

By choosing one of the available alternatives, a new regexp that defines a smaller language than the original one, is obtained.

$$e_k, 1 \le k \le n, \qquad \text{choice for the union} \quad e_1 \bigcup ... \bigcup e_k \bigcup ... e_n$$

$$e \qquad \text{choice for the iteration} \quad e^*, e^+, e^n$$

$$\varepsilon \qquad \text{choice for the iteration} \quad e^*$$

Given a regexp $e_1$, one can always DERIVE another regexp $e_2$ by replacing a subexpression (short s.e.) of the regxexp $e_1$ with one of the possible choices.

DERIVATION RELATION between two regexps *e'* and *e"*

$$e' \Rightarrow e''$$ If the two regexp *e'* and *e"* can be factored as

$$e' = \alpha\beta\gamma \qquad e'' = \alpha\delta\gamma$$

with $\beta$ s.e. of $e'$, $\delta$ s.e. of $e''$, $\delta$ is a choice of $\beta$

DERIVATION STEPS can be applied in sequence, two or more times:

$$e_0 \overset{n}{\Rightarrow} e_n \quad \text{if} \quad e_0 \Rightarrow e_1, e_1 \Rightarrow e_2, \;\;\ldots\;, \; e_{n-1} \Rightarrow e_n$$

$$e_0 \overset{+}{\Rightarrow} e_n \quad e_0 \;\; \text{derives} \;\; e_n \;\; \text{in } n \geq 1 \quad \text{steps}$$

$$e_0 \overset{*}{\Rightarrow} e_n \quad e_0 \;\; \text{derives} \;\; e_n \;\; \text{in } n \geq 0 \quad \text{steps}$$

VARIOUS EXAMPLES

One-step derivation and multiple-step derivation:

$$a^* \bigcup b^+ \Rightarrow a^* \qquad a^* \bigcup b^+ \Rightarrow b^+$$

$$a^* \bigcup b^+ \Rightarrow a^* \Rightarrow \varepsilon \qquad a^* \bigcup b^+ \overset{2}{\Rightarrow} \varepsilon \qquad a^* \bigcup b^+ \overset{+}{\Rightarrow} \varepsilon$$

$$a^* \bigcup b^+ \Rightarrow b^+ \Rightarrow bbb \qquad a^* \bigcup b^+ \overset{2}{\Rightarrow} bbb \qquad a^* \bigcup b^+ \overset{+}{\Rightarrow} bbb$$

Of the regexp that are obtained by derivation from the initial regexp $e$,
some may contain letters of the alphabet and also metasymbols,
some may contain only letters of the alphabet. The latter constitute
the language generated by the initial regexp $e$.

THE LANGUAGE $L(r)$ DEFINED
(or GENERATED) BY A GIVEN REGEXP $r$ IS:

$$L(r) = \left\{ x \in \Sigma^* \mid r \overset{*}{\Rightarrow} x \right\}$$

Two regexps are EQUIVALENT if they define (generate) the same language.

THE LANGUAGE DEFINED BY A DERIVED REGEXP IS CONTAINED
IN THE LANGUAGE DEFINED BY THE DERIVING REGEXP

There may exist several derivation orders that generate the same
string, which are however substantially equivalent.

EXAMPLES:

$1. (ab)^* \Rightarrow abab$

$2. (ab \cup c) \Rightarrow ab$

$3. a(ba \cup c)^* d \Rightarrow ad$

$4. a(ba \cup c)^* d \Rightarrow a(ba \cup c)(ba \cup c)d$

$5. a^*(b \cup c \cup d) f^+ \Rightarrow aaa(b \cup c \cup d) f^+$

$6. a^*(b \cup c \cup d) f^+ \Rightarrow a^* c f^+$

$7. a^*(b \cup c \cup d) f^+ \overset{+}{\Rightarrow} aaacf^+$   in 2 steps

$8. a^*(b \cup c \cup d) f^+ \overset{+}{\Rightarrow} aaacff$   in 3 steps

## AMBIGUITY OF REGULAR EXPRESSIONS

A string (phrase) may be obtained by two derivations, that differ not only in the ordere of choices, but in a more substantial way.

Examples:

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

$$(a \cup b)^* a (a \cup b)^* \Rightarrow (a \cup b) a (a \cup b)^* \Rightarrow aa(a \cup b)^* \Rightarrow aa\varepsilon \Rightarrow aa$$

$$(a \cup b)^* a (a \cup b)^* \Rightarrow \varepsilon a (a \cup b)^* \Rightarrow \varepsilon a (a \cup b) \Rightarrow \varepsilon aa \Rightarrow aa$$

A regexp $f$ is AMBIGUOUS if and only if, the corresponding marked regexp $f'$ generates two marked strings $x$ and $y$ such that, if the indices are removed, the underlying ordinary strings are identical.

$$f' = (a_1 \cup b_2)^* a_3 (a_4 \cup b_5)^*$$

defines a regular language over the (marked) alphabet

$$\{a_1, b_2, a_3, a_4, b_5\}$$

$a_1 a_3$ , $a_3 a_4$ exhibit the ambiguity phenomenon

EXAMPLE (ambiguity)

$(aa \mid ba)^* a \mid b(aa \mid ba)^*$    is ambiguous

 in fact, the marking

$(a_1 a_2 \mid b_3 a_4)^* a_5 \mid b_6 (a_7 a_8 \mid b_9 a_{10})^*$

 generates the two strings

$b_3 a_4 a_5$     $b_6 a_7 a_8$

 that project ambiguously onto the same phrase:     $baa$

## OTHER OPERATORS

Basic operators: union, concatenation, star.
Derived operators: power, cross.

$$e^h = ee...e \quad\quad e^+ = ee^*$$
$$\underset{h \text{ times}}{}$$

EXAMPLE: floating point fractional numbers, with or without sign and exponent

$$\Sigma = \{+,-,\bullet,E,d\}$$

$$r = s.c.e$$

$$s = (+\cup-\cup\varepsilon) \text{ sign } \pm, \text{ optional}$$

$$c = (d^+ \bullet d^* \cup d^* \bullet d^+) \text{ constant, integer or fractional}$$

$$e = (\varepsilon \cup E(+\cup-\cup\varepsilon)d^+) \text{ exponent, optional, preceded by E}$$

$$(+\cup-\cup\varepsilon)(d^+ \bullet d^* \cup d^* \bullet d^+)(\varepsilon \cup E(+\cup-\cup\varepsilon)d^+)$$

$$+dd \bullet E - ddd \quad\quad +12 \bullet E - 341 \quad 12.10^{-341}$$

REPETITION: from $k$ to $n > k$ times
OPTIONALITY:
INTERVAL OF AN ORDERED SET:

$$[a]_k^n = a^k \bigcup a^{k+1} \bigcup ... \bigcup a^n$$

$$[a] = (\varepsilon \bigcup a)$$

$$(0...9) \quad (a...z) \quad (A...Z)$$

By admitting in a regexp also the presence of the set-theoretic operations INTERSECTON, COMPLEMENT and DIFFERENCE, one obtains the so-called EXTENDED REGULAR EXPRESSIONS.

EXAMPLE (intersection) – allows to express the request that the phrases of the language satisfy two conditions (both, not either one).

the same result of intersection but obtained in a more complicated way
(*bb* is sourrounded by two strings both of either even or odd length)

↓

| | |
|---|---|
| alphabet | $\{a, b\}$ |
| string contains $bb$ | $(a\,|\,b)^* bb(a\,|\,b)^*$ |
| string length is even | $((a\,|\,b)^2)^*$ |

$$(a\,|\,b)^* bb(a\,|\,b)^* \bigcap ((a\,|\,b)^2)^*$$

$$((a\,|\,b)^2)^* bb((a\,|\,b)^2)^* \,|\, (a\,|\,b)((a\,|\,b)^2)^* bb(a\,|\,b)((a\,|\,b)^2)^*$$

EXMPLE (complement and intersection)

The regexp $r = (\,a\,b\,)^*$ generates the strings that
- • do not start by $b$
- • do not end by $a$
- • do not contain either the substring $aa$ or the substring $bb$

$$r' = \neg(b(a \cup b)^* \cup (a \cup b)^* a \cup (a \cup b)^*(aa \cup bb)(a \cup b)^*)$$

De Morgan theorem:

$$r' = \neg b(a \cup b)^* \cap \neg(a \cup b)^* a \cap \neg((a \cup b)^*(aa \cup bb)(a \cup b)^*)$$

EXAMPLE: how to transform an extended regexp over the alphabet { $a$, $b$, $c$ }
into an equivalent regexp, such that it uses only the basic operators U, $*$ and .

$$\neg((a \mid b)(\neg\varnothing))$$
$$= \neg((a \mid b)(a \mid b \mid c)^*) = \varepsilon \mid c(a \mid b \mid c)^*$$

# CLOSURE OF THE REG FAMILY WITH RESPECT TO OPERATIONS / 1

Let θ be an operator, that applied to one language (operand), or to a pair of languages, produces another language (result).

A FAMILY OF LANGUAGES IS SAID TO BE CLOSED WITH RESPECT TO AN OPERATOR θ, IF THE RESULT LANGUAGE BELONGS TO THE SAME FAMILY AS THE OPERAND LANGUAGE(S).

---

PROPERTY: the family REG of the regular languages is CLOSED with respect to the following operators: concatenation, union, star, and therefore is closed also w.r.t. the derived operators, e.g. cross, repetition, optionality, etc (this follows from the definition of regexp itself).

---

This implies that any two regular languages $L_1$ and $L_2$ can be combined by means of the above mentioned operators, and yet the resulting language remains in the family of the regular languages.

AN EVEN STRONGER PROPERTY: the family REG of the regular languages is the *smallest* family of languages such that both the following properties hold:

- REG contains all finite languages
- REG is closed w.r.t. concatenation, union and star

Later the LIB family will be defined (of the so-called context-free languages), which is closed w.r.t. concatenation, union and star as well, but is not the smallest such family, because the following strict containment holds:

$$REG \subset LIB$$

Moreover, REG is closed w.r.t. INTERSECTION, COMPLEMENT (and DIFFERENCE) and MIRRORING (this is difficult to prove directly, but is esay by using the concepts and tools of automata theory).

## LINGUISTIC ABSTRACTION

Linguistic abstraction transforms the phrases of a real, effective language, and gives them a simpler form, called abstract representation.

To do so, the symbols of the effective alphabet are discarded and replaced by those of the abstract alphabet.

AT THE ABSTRACT LEVEL, THE TYPICAL STRUCTURES OF MOST ARTIFICIAL LANGUAGES CAN BE OBTAINED BY THE COMPOSITON OF FEW ELEMENTARY PARADIGMS, AND BY THE BASIC LANGUAGE OPERATIONS SUCH AS CONCATENATION, UNION AND ITERATION.

From the abstract language to the effective language → choice of the actual lexical elements (e.g. keywords, identifiers, ...)

COMPILER DESIGN MAKES REFERENCE TO THE ABSTRACT LANGUAGE TO PROCESS, RATHER THAN TO THE EFFECTIVE ONE.

Artificial languages (e.g. programming languages) contain few abstract structures, among which lists play a relevant role; lists can be esaily modeled by regexps.

A list is a sequence of a number of elements, not fixed in advance.

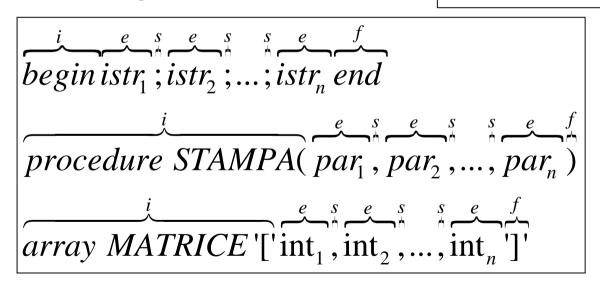A list can be generated by the regexp $e^+$ or $e^*$, if the empty list is admitted.

The element $e$ can be a teminal ($=$ alphabetical symbol ) or a compound object of some kind (for instance, a list itself …).

LISTS WITH SEPARATORS AND START-MARKER OR END-MARKER

EXAMPLES:

$$ie(se)^* f \qquad i[e(se)^*]f$$

$$\underbrace{begin}_{i}\ \underbrace{istr_1}_{e}\ \underbrace{;}_{s}\ \underbrace{istr_2}_{e}\ \underbrace{;}_{s} \dots \underbrace{;}_{s}\ \underbrace{istr_n}_{e}\ \underbrace{end}_{f}$$

$$\underbrace{procedure\ STAMPA(}_{i}\ \underbrace{par_1}_{e}\ \underbrace{,}_{s}\ \underbrace{par_2}_{e}\ \underbrace{,}_{s} \dots \underbrace{,}_{s}\ \underbrace{par_n}_{e}\ \underbrace{)}_{f}$$

$$\underbrace{array\ MATRICE\ '['}_{i}\ \underbrace{int_1}_{e}\ \underbrace{,}_{s}\ \underbrace{int_2}_{e}\ \underbrace{,}_{s} \dots \underbrace{,}_{s}\ \underbrace{int_n}_{e}\ \underbrace{']'}_{f}$$

SUBSTITUTION (REPLACEMENT): operation that replaces the terminal
characters of the source language with the phrases of the target
(or destination) language.

$$L \subseteq \Sigma^* \quad L_b \subseteq \Delta^*$$

$$x \in L \quad x = a_1 a_2 ... a_n \quad \text{and for some} \quad a_i = b$$

Substituting the language $L_b$ to the letter $b$ in the string $x$ produces
a language over the alphabet $(\Sigma \setminus \{ b \}) \cup \Delta$, defined as follows:

$$\{ y \mid y = y_1 y_2 ... y_n \wedge (\text{if } a_i \neq b \text{ then } y_i = a_i \text{ else } y_i \in L_b \}$$

NESTED LISTS (sometimes called PRECEDENCE LISTS)

Lists may contain atomic objects
as elements, but also other lists,
of lower level.

$$lista_1 = i_1 lista_2 (s_1 lista_2)^* f_1$$

$$lista_2 = i_2 lista_3 (s_2 lista_3)^* f_2$$

$$...$$

$$lista_k = i_k e_k (s_k e_k)^* f_k$$

EXAMPLES

livello 1: $begin \; istr_1 ; istr_2 ; ... ; istr_n \; end$

livello 2: $STAMPA(var_1, var_2, ..., var_n)$

$$3 + \underbrace{5 \times 7 \times 4}_{monomio1} - \underbrace{8 \times 2 \div 5}_{monomio2} + 8 + 3$$

padre, madre, figlio e figlia
un padre forte, severo e giusto, una madre amorevole e fedele

un libro come lista di capitoli separati da pagine bianche,
                                        chiusa tra due copertine
un capitolo come lista di sezioni

# Bibliography

– S. Crespi Reghizzi, *Linguaggi Formali e Compilazione*, Pitagora Editrice Bologna, 2006

– Hopcroft, Ullman, *Formal Languages and their Relation to Automata*, Addison Wesley, 1969

– A. Salomaa – *Formal Languages*, Academic Press, 1973

– D. Mandrioli, C. Ghezzi – *Theoretical Foundations of Computer Science*, John Wiley & Sons, 1987

– L. Breveglieri, S. Crespi Reghizzi, *Linguaggi Formali e Compilatori: Temi d'Esame Risolti,* web site (eng + ita)