

* Today: Regular expressions

Fix an alphabet Σ

* Def: A regular expression (regex) r a string in the letters of Σ , together with the symbols " ϵ ", " $|$ ", " $*$ ", " ϕ ", [and " $($ ", " $)$ "] satisfying one of the following:

- (1) $r = \phi$
- (2) $r = \epsilon$
- (3) $r = a$ for some $a \in \Sigma$
- (4) $r = r_1 r_2$ where r_1, r_2 are also regexes
- (5) $r = r_1 | r_2$ where r_1, r_2 are regexes
- (6) $r = r_1^*$

[In any of these options, $()$ signify grouping]
 [(7) $r = (r_1)$ where r_1 is a regular expression]

Just like in an algebraic expressions-

** We assume that " $|$ ", " $*$ ", " ϕ ", " $($ ", " $)$ " are not in Σ .

** Order of operations:

Brackets first, then $*$, then concatenation, then " $|$ " (or)

* Examples. Let $\Sigma = \{0, 1\}$

$r = \phi, r = \epsilon, r = 0, r = 1$

$r = \phi^*, r = 0^*, r = \epsilon^*$

$r = 0|1, r = 0|1|0^* \# \left[\begin{array}{c} (0|1)|0^* \# 0|(1|0^*) \\ \uparrow \qquad \qquad \uparrow \\ \text{equivalent to } 0|1|0^* \end{array} \right]$

$r = (0|1|0^*)^* 00|1$

$r = \underbrace{(0|1|\phi^*|1|0^*)}_{r_1}^* \underbrace{0|10}_{r_2} \underbrace{(0|1|\epsilon)}_{r_2}$

$r = \underbrace{r_1^*}_{r_1} \underbrace{r_2}_{r_2}$

$r_1 = \underbrace{(0|1)}_{r_3} | \underbrace{\phi^*}_{r_4} | \underbrace{1|0^*}_{r_5} = r_3 | r_4^* | r_5$

[Continue breaking up the expression mentally until you hit either ϕ, ϵ , or a letter.]

** Matching

Let r be a regex. Let $w \in \Sigma^*$ be a string.

We say that w matches r if one or more of the following hold:

- (1) $r = \varepsilon$ and $w = \varepsilon$
- (2) $r = a$ for some $a \in \Sigma$, and $w = a$
- (3) $r = r_1 r_2$ for regexes r_1, r_2 , and $w = xy$, where x, y are strings, and x matches r_1 , and y matches r_2 .
- (4) $r = r_1 | r_2$ and w either matches r_1 or r_2 (or both).
- (5) $r = r_1^*$, and either $w = \varepsilon$ or $w = x_1 x_2 \dots x_k$ where each x_i is a string, and each x_i matches r_1 .

* No string matches $r = \phi$.

** Examples

- $r = 0$: $w = 0$ only string that matches
- $r = 1$: $w = 1$ " " " "
- $r = \varepsilon$: $w = \varepsilon$ " " " "

$r = 010$: $w = 010$ only match

$r = 1\phi$: nothing matches!

③

$r = 0|1$: $w = 0, w = 1$ only matches

$r = 1^*$: $w = \varepsilon, w = 1, w = 11, w = 111, \dots$

$r = (01)^*$: $\varepsilon, 01, 0101, 010101, \dots$

$r = (00|11)^*$: $\varepsilon, 00, 11, 0000, 1111, \underline{0011}, 1100, \dots$

$w = \underline{00}1\underline{0}$ does not match.

~~$r = 01^* | 10$~~

$r = 01^* | 0^*1$: $w = 001, 011, 01 \checkmark$

$w = 0101$ not a match

$w = 1 \checkmark$

** The language of a regex

Let r be a regex.

The language of r , denoted $L(r)$ is the set of all strings that match r .

E.g. $r = 0$, $L(r) = \{0\}$

$r = \phi$, $L(r) = \phi$

$r = \varepsilon$, $L(r) = \{\varepsilon\}$

$r = 0(1|0)^*1$, $L(r) =$ strings starting with 0 and ending with 1.

④