

* Today : Regular expressions

Fix an alphabet Σ'

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

$$(1) \quad r = \phi$$

$$(2) \quad r = \epsilon$$

$$(3) \quad r = a \text{ for some } a \in \Sigma$$

$$(4) \quad r = r_1 r_2 \text{ where } r_1, r_2 \text{ are also regexes}$$

$$(5) \quad r = r_1 | r_2 \text{ where } r_1, r_2 \text{ are regexes}$$

$$(6) \quad r = r_1^*$$

[In any of these options, () signify grouping]

[$(7) \quad r = (r_1)$ 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 $\Sigma = \{0, 1\}$

(2)

$$r = \phi, \quad r = \varepsilon, \quad r = 0, \quad r = 1$$

$$r = \phi^*, \quad r = 0^*, \quad r = \varepsilon^*$$

$$r = 0|1, \quad r = 0|1|0^* \quad \boxed{(0|1)|0^* \neq 0|(1|0^*)}$$

↑ ↑
equivalent to $0|1|0^*$

$$r = (010|1)^* 00|1$$

$$r = \underbrace{(01|\phi^*|110)}_{r_1}^* \underbrace{010}_{\substack{\downarrow \\ r_2}} \underbrace{(011|\varepsilon)}_{\substack{\downarrow \\ r_2}}$$

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

$$r_1 = \underbrace{01}_{r_3} \mid \underbrace{\phi^*}_{r_4} \mid \underbrace{110}_{r_5} = r_3 \mid r_4^* \mid r_5$$

[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) $r = \epsilon$ and $w = \epsilon$
- (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 $w = \epsilon$ 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 $r = \emptyset$.

**** Examples**

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

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

- $r = 1\emptyset$: nothing matches!

- $r = 0 \mid 1$: $w = 0, w = 1$ only matches
- $r = 1^*$: $w = \epsilon, w = 1, w = 11, w = 111, \text{etc}$
- $r = (01)^*$: $\epsilon, 01, 0101, 010101, \dots$
- $r = (00|11)^*$: $\epsilon, 00, 11, 0000, 1111, \underline{0011}, 1100, \dots$
 $w = \underline{0010}$ does not match.

~~$r = 01^* \mid 10$~~

- $r = 01^* \mid 0^* 1$ $w = 001, 011, 01 \checkmark$
 $w = 0101$ 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. $r = 0$, $L(r) = \{0\}$

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

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

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