# Linguagens e Programação

Ana Madureira amd@isep.ipp.pt António Silva ass@isep.ipp.pt
Bruno Cunha cun@isep.ipp.pt José Marinhojsm@isep.ipp.pt
Paulo Ferreira pdf@isep.ipp.pt

Instituto Superior de Engenharia do Porto

2021/2022

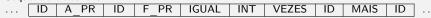
#### Análise Léxica

- Agrupar os caracteres individuais do programa fonte em "palavras", designadas por tokens, que são as unidades elementares das linguagens de programação, e que serão usados pelo analisador sintáctico (parser)
- A sequência de caracteres que formam um token chama-se lexema
- Eliminar sequências de caracteres inúteis, que não têm significado sintático (espaços, tabulações, LF, CR) e os comentários
- produzir mensagens de erro quando algum carácter ou sequência de caracteres não são reconhecidos
- localizar os tokens no texto fonte (nº de linha, nº de coluna)
- actuar como pré-processador (em certos compiladores)

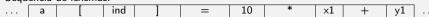
#### Exemplo



Sequência de tokens:



Sequência de lexemas:



#### Tokens e Lexemas

Token - representa um conjunto de cadeias de entrada possível

Lexema – é uma determinada cadeia de entrada associada a um token

Tokens	Lexemas
FOR	for
IF	if
WHILE	while
NUMBER	1089, 142857, 0.2, 3.14159
IDENTIFIER	i, j, counter, StudentName
OP PLUS	+
OP GREATER EQUAL	>=
OPEN_PAR	(

#### Reconhecimento dos Tokens

#### Expressões Regulares — representam padrões de cadeias de caracteres

- Os tokens são geralmente especificados utilizando expressões regulares
- As expressões regulares são regras bem definidas que geram um conjunto de sequências de caracteres (strings)
- O conjunto de strings gerado por uma expressão regular chama-se uma linguagem regular

**Autómatos Finitos** — forma matemática de descrever tipos particulares de algoritmos, para fazer o reconhecimento de padrões em cadeias de caracteres

Alfabeto: Qualquer conjunto finito, não vazio, de símbolos.

- ullet Um alfabeto é designado por  $\Sigma$
- Considerem-se os seguintes exemplos:

$$\Sigma = \{0,1\}$$
 , alfabeto para numeros binários.

$$\Sigma = \{a, b, c, \dots, z\}$$
, o conjunto das letras minúsculas.

Palavra ou cadeia: é uma sequência finita (eventualmente vazia) de símbolos do alfabeto  $\Sigma$ 

- A palavra vazia (  $\varepsilon$ ) contém zero ocorrências de símbolos
- |W| designa o comprimento da palavra W.

$$|aa| = 2$$

$$|\varepsilon| = 0$$

**Potência de um alfabeto:**  $\Sigma^k$  representa o conjunto de todas as palavras de comprimento k, sobre um alfabeto  $\Sigma$ 

```
Por exemplo, para: \Sigma=\{0,1\}: \Sigma^0=\{\varepsilon\} \mbox{ (cont\'em apenas a palavra vazia)} \Sigma^1=\{0,1\} \mbox{ (n\~ao confundir com $\Sigma$)} \Sigma^2=\{00,01,10,11\} \Sigma^3=\{000,001,010,011,100,101,110,111\}
```

**Fecho de um alfabeto**  $\Sigma$ : representado por  $\Sigma^* = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \Sigma^3 \cup ...$ , corresponde ao conjunto de todas as palavras que podem ser formadas com os símbolos de  $\Sigma$ , incluindo a palavra vazia  $\varepsilon$ ;

Para o alfabeto 
$$\Sigma = \{0, 1\}$$
:

$$\Sigma^* = \left\{0,1\right\}^* = \left\{\varepsilon,0,1,00,01,10,11,000,001,\ldots\right\}$$

Concatenação de palavras: Justapondo duas palavras  $\alpha$  e  $\beta$  será obtido  $\alpha\beta$ , uma palavra formada pelos símbolos de  $\alpha$  seguidos dos símbolos de  $\beta$ .

A operação de concatenação de palavras possui as seguintes propriedades:

• A palavra vazia  $\varepsilon$  constitui o elemento identidade (neutro) da operação, isto é, para qualquer palavra  $\alpha$ :

$$\alpha \varepsilon = \varepsilon \alpha = \alpha$$

• **é associativa**, ou seja, para quaisquer palavras  $\alpha$ ,  $\beta$  e  $\gamma$ :

$$(\alpha\beta)\gamma = \alpha(\beta\gamma)$$

• **não é comutativa**, ou seja,  $\alpha\beta \neq \beta\alpha$ , sempre que  $\alpha \neq \beta$ .

$$lpha=$$
 aa  $eta=$  b  $lphaeta=$  aab  $eq$  baa  $=$   $etalpha$ 

- O comprimento é **aditivo** relativamente à concatenação,  $|\alpha\beta| = |\alpha| + |\beta|$
- Duas palavras são iguais sempre que uma for a cópia, símbolo a símbolo, da outra;
- Dadas duas palavras  $\alpha$  e  $\beta$ :
  - ightharpoonup a palavra  $\alpha$  é **prefixo** da palavra  $\alpha\beta$ ,
  - **a** palavra  $\beta$  é **sufixo** da palavra  $\alpha\beta$ ;
- Diz-se que a palavra  $\alpha$  é parte (**subpalavra** ou **factor**) da palavra  $\beta$  sempre que existam as palavras  $\gamma$  e  $\delta$  , tais que  $\beta=\gamma\alpha\delta$ ;

• Para qualquer número natural n e qualquer palavra  $\alpha$ , designa-se por  $\alpha^n$  a **potência** de ordem n

Por exemplo, para o alfabeto  $\Sigma = \{a, b, c, \dots, z\}$  e  $\alpha = ab$  teremos:

$$lpha^1 = ab$$
 $lpha^2 = abab$ 
 $lpha^3 = ababab$ 

Note-se que  $(ab)^3=ababab$ , enquanto que  $ab^3=abbb$ . Por convenção para qualquer palavra  $\alpha$  temos que  $\alpha^0=\varepsilon$ ;

• Designa-se por  $\alpha^R$  a **palavra inversa** de  $\alpha$ , ou seja, se  $\alpha = ab$  então  $\alpha^R = ba$ 

$$(\alpha^R)^R = \alpha$$
  
 $(\alpha^R)^n = (\alpha^n)^R$  para qualquer  $n$ 

- Uma linguagem formal é um conjunto de cadeias de símbolos (palavras) de um determinado alfabeto
- Uma linguagem é finita se as suas frases formam um conjunto finito (caso contrário, a linguagem é infinita)
- Uma linguagem infinita precisa ser definida através de uma representação finita

Por exemplo, a linguagem dos números naturais menores que 10 é finita, e pode ser representada, literalmente, como:

$$L = \{1, 2, 3, 4, 5, 6, 7, 8, 9\}$$

Já a linguagem dos números naturais como um todo não é uma linguagem finita, já que existem infinitos números naturais

Uma linguagem L sobre um alfabeto  $\Sigma$  é um conjunto de palavras de  $\Sigma^*$ , isto é, é um subset de  $\Sigma^*$  ( $L \subseteq \Sigma^*$ )

Uma linguagem L diz-se que:

- L pode não incluir todas as palavras de  $\Sigma^*$ , nem sequer as suas palavras incluírem todos os símbolos de  $\Sigma$ ;
- Uma **linguagem vazia**, que não contém palavras, designa-se por ∅.

L=
$$\emptyset$$
 é diferente de  $L = \{\varepsilon\}$ 

- A linguagem designa-se como **completa** se coincide com a totalidade do conjunto  $\Sigma^*$ ;
- Se L e M são linguagens de alfabeto Σ, a concatenação de L com M é uma linguagem de alfabeto Σ, definida por

$$LM = \{\alpha\beta \mid \alpha \in L \text{ e } \beta \in M\}$$
 
$$\{a, ab\} \{b, ba\} = \{ab, aba, abb, abba\}$$

• As **potências**  $L^n$  de uma linguagem L são definidas indutivamente por:

$$L^{0} = \{\varepsilon\}$$

$$L^{n+1} = LL^{n}, \text{para } n \ge 1$$

$$L^{n} = \{x_{1}x_{2} \dots x_{n} | x_{i} \in L, 1 \le i \le n\}$$

Isto é, a linguagem das palavras obtidas por concatenação de  $\it n$  palavras de  $\it L$ .

 O fecho de Kleene de uma linguagem L é a reunião de todas as potências finitas de L.

$$L^* = L^0 \cup L^1 \cup L^2 \cup L^3 \cup \ldots = \cup_{n \geq 0} L^n$$

é conjunto das palavras que se podem formar por concatenação de zero ou mais palavras de  ${\it L}$ 

 $\Sigma^*$ , não é mais do que o fecho de *Kleene* de  $\Sigma$ .

- O fecho positivo da linguagem L, definido por  $L^+ = L^1 \cup L^2 \cup L^3 \cup ...$ 
  - Se  $\varepsilon \in L$  então  $L^+ = L^*$
  - ▶ Se  $\varepsilon \notin L$  então  $L^+ = L^* \setminus \{\varepsilon\}$
- Toda a expressão regular representa uma linguagem regular;
- Toda a linguagem regular é representável por alguma expressão regular.

Para o alfabeto binário  $\Sigma = \{0, 1\}$ , os seguintes sets são linguagens:

- ▶ ∅ linguagem vazia;
- $\{\varepsilon\}$  linguagem constituída pela palavra nula;
- ▶ {0} linguagem constituída por uma única palavra "0";
- $ightharpoonup \{00,01,10,11\}$  linguagem constituída por palavras de comprimento 2.

### Equivalência entre Expressões Regulares e Linguagens

O conjunto de Expressões Regulares sobre um alfabeto  $\Sigma$  e o conjunto das linguagens por elas descritas são definidos indutivamente por:

- $\varepsilon$  é uma expressão regular sobre o alfabeto  $\Sigma$  e descreve a linguagem  $\{\varepsilon\}$ ,  $L(\varepsilon) = \{\varepsilon\}$ ;
- $\emptyset$  é uma expressão regular sobre o alfabeto  $\Sigma$  e descreve a linguagem vazia  $\emptyset$  isto é,  $L(\emptyset) = \emptyset$ ;
- Se  $a \in \Sigma$  então a é uma expressão regular sobre  $\Sigma$  e descreve a linguagem  $\{a\}$ , isto é,  $L(a) = \{a\}$ ;

## Operações sobre ERs/LF

• Se r e s são Expressões Regulares sobre  $\Sigma$  que descrevem as linguagens L(r) e L(s):

União 
$$\longrightarrow r|s=L(r)\cup L(s)$$
  
Concatenação  $\longrightarrow rs=L(r)L(s)$   
Fecho de Kleene  $\longrightarrow r^*=L(r)^*$ 

Se r e s que permitem reconhecer as palavras,  $\{aa, ab\}$  e  $\{bb, ba\}$ :

#### União

$$r|s = L(r) \cup L(s) = \{aa, ab, bb, ba\}$$

#### Concatenação

$$rs = L(r)L(s) = \{aabb, aaba, abbb, abba\}$$

#### Fecho de Kleene

$$r^* = L(r)^* = \{ \varepsilon, aa, ab, aaaa, aaab, abaa, abab, aaaaaaa, aaaaab, \ldots \}$$

#### Expressões Regulares/Linguagens Formais

- Qualquer expressão regular r representa uma linguagem que se designa por L(r);
- Diz-se que duas Expressões Regulares são equivalentes se, e só se, as linguagens por elas descritas são iguais.

### Expressões Regulares — Precedências

A precedência para as operações união (|), concatenação (·) e fecho de Kleene (\*) correspondem às usadas nas expressões aritméticas para a **adição**, a **multiplicação** e a **potenciação**, respectivamente.

#### fecho de Kleene > concatenação > união.

Considerando as "ERs" r, s e t. Então a "ER"  $rs^*|t$  será avaliada da forma seguinte:

- 1º Passo avaliação da potência (fecho de Kleene)  $rs^*|t\Rightarrow r\{\varepsilon,s,ss,sss,\ldots\}|t$
- 2º Passo avaliação da concatenação  $rs^*|t\Rightarrow \{r\varepsilon, rs, rss, rsss, \ldots\}|t$  sendo que  $r\varepsilon=r$
- 3° Passo avaliação da união  $rs^*|t\Rightarrow \{r,rs,rss,rsss,\ldots,t\}$

## Exemplos de Expressões Regulares

Seja $\Sigma = \{0,1\}$	
expressão regular	Linguagem Descrita
0 1	{0,1}
0*	$\{\varepsilon, 0, 00, 000, 0000, 00000, \ldots\}$
0*11	$\{11,011,0011,00011,000011,0000011,\ldots\}$
(01)*	$\{\varepsilon, 01, 0101, 010101, 01010101, 0101010101$
$(0 1)^*$	$\{\varepsilon, 0, 1, 00, 01, 10, 11, 000, 001, 010, 011, 100, \ldots\}$
1*0(0 1)*	$\{0, 10, 01, 00, 000, 100, 010, 001, 110, 101, 011, \ldots\}$
$(0 1)^*0$	sequências de palavras terminadas (ou com sufixo) em $0$
$(0 1)^*1$	sequências de palavras terminadas (ou com sufixo) em $1$

sequência que contém o factor ou subpalavra 101

#### Propriedades das Expressões Regulares

#### Sejam u, v e x Expressões Regulares. Tem-se que:

$$u|v = v|u$$

$$u|u = u$$

$$(u|v)|x = u|(v|x)$$

$$u\varepsilon = \varepsilon u = u$$

$$(uv)x = u(vx)$$

$$u(v|x) = uv|ux$$

$$(u|v)x = ux|vx$$

$$\varepsilon^* = \varepsilon$$

## Propriedades das Expressões Regulares

Sejam u, v e x Expressões Regulares. Tem-se que:

$$u^* = u^*u^* = (u^*)^* = u^*|u^*$$

$$u^* = \varepsilon|u^* = (\varepsilon|u)^* = (\varepsilon|u)u^* = \varepsilon|uu^*$$

$$u^* = (u|\dots|u^k)^*, k \ge 1$$

$$u^* = \varepsilon|u|\dots|u^{k-1}|u^ku^*, k > 1$$

$$u^*u = uu^*$$

$$(u|v)^* = (u^*|v^*)^* = (u^*v^*)^* = (u^*v)^*u^* = u^*(vu^*)^*$$

$$u(vu)^* = (uv)^*u$$

$$(u^*v)^* = \varepsilon|(u|v)^*v$$

#### Expressões Regulares — Aplicação prática FLEX

O carácter "x" X Qualquer carácter excepto mudança de linha nmudança de linha [xyz] Um dos caracteres "x", "y" ou "z" A cadeia de caracteres "xyz" xyz Um dos caracteres no intervalo de "a" a "z" ou de "A" a "Z" [a-zA-Z]Qualquer um dos operadores "-", "+", "\*" ou "/", sendo que o [-+\*/] símbolo "-" tem de aparecer em primeiro lugar dada a possibilidade de ambiguidade com a definição de intervalo Um dos caracteres "a", "b" ou de "i" a "o" ou "Z" [abj-oZ]  $[^A-Z n]$ Qualquer carácter excepto no intervalo de "A" a "Z" ou mudança de linha O carácter "r" zero ou mais vezes r\* O carácter "r" uma ou mais vezes r+ O carácter "r" zero ou uma vez r?

## Expressões Regulares — Aplicação prática FLEX

r{2,5} r{2,}	O carácter "r" repetido de duas a cinco vezes O carácter "r" repetido pelo menos duas vezes
r{4}	O carácter "r" repetido exactamente quatro vezes
{macro}	Substituição/Expansão da macro definida anteriormente
(r)	O carácter "r", sendo que os parêntesis permitem estipular precedências
xyz*	A sequência "xy" seguida de zero ou mais "z"s
(xyz)*	A sequência "xyz" repetida zero ou mais vezes
r s	O carácter "r" ou "s" (alternativa)
^r	O carácter "r" apenas se no início da linha
r\$	O carácter "r" apenas se no final da linha (não consome o \n)
^xyz\$	Uma linha que contém apenas a cadeia de caracteres "xyz"
< <e0f>&gt;</e0f>	Fim de ficheiro

## Expressões Regulares — Exemplos

- Definição de um identificador
  - $\blacktriangleright \ \mathsf{letra} \to \mathsf{[A-Za-z]}$
  - ▶ dígito  $\rightarrow$  [0-9]
  - ▶ identificador  $\rightarrow$  {letra} ({letra}|{digito})\*
- Definição de um número com parte decimal
  - ▶ dígito → [0-9]
  - ▶ dígitos → {dígito}+
  - ▶ parteFraccionária → (\. {dígitos})?
  - ightharpoonup expoente ightharpoonup (E [+-]? {dígitos})?
  - ▶ num → {dígitos} {parteFraccionária} {expoente}

### Expressões Regulares — Exercícios

- Escreva uma expressão regular para cada um dos seguintes casos, usando  $\Sigma = \{0,1\}$ :
  - Sequência com sufixo 1
  - Sequência com prefixo 1
  - Sequência que contém dois 1s
  - Sequência que contém pelos menos três 1s e um 0
  - Sequência que contém o factor 01
  - Sequência que começa e termina com o mesmo algarismo