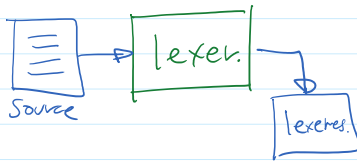
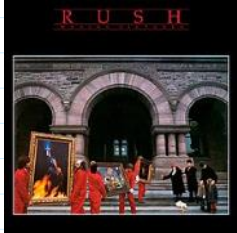


4 Regular Expressions

Wednesday, August 30, 2023 10:56 AM



Linguistics:-

- Orthography:- How to write words?
apple manzana pomme U Y O while
for int
- Syntax:- How to write sentences?
"has jump apple dog green"
- Semantics:- What do sentences mean?
"colorless green ideas sleep furiously"
"the chickens are ready to eat"
"moving pictures" 

• How to WRITE WORDS?

- FOR
- WHILE
- IF
- user defined names: "identifiers"
- "Literals"
 - numerical values
 - integers
 - real numbers
 - complex numbers
 - scientific numbers
 - strings
 - arrays
 - dictionaries
 - tuples

• MATHEMATICAL LANGUAGES.

- S. Koenig
- E. Post
- Alphabet:- finite.
a set of symbols.
- Language:- a set of sequences/words made from symbols in an alphabet.

e.g. $\Sigma = \{a, b\}$

$L = \{a, aa, aba, bbaa\}$

= Specification problem
How to precisely describe a Language

$\Sigma = \{a, b\}$

$L = \{a, ab, bab, abb, \dots\}$

= Recognition Problem.

given a word w and Language L ,
 $w \in L?$

SPECIFICATION: REGULAR EXPRESSIONS

$\Sigma, \mid, *, ()$

Regular Expressions:

- a word w is a regular expression. $\{w\}$
- if r_1 and r_2 are regular expressions
 - $r_1 r_2$ $\{w_1 w_2 \mid w_1 \in r_1 \wedge w_2 \in r_2\}$
 - $r_1 \mid r_2$ $\{w \mid w \in r_1 \vee w \in r_2\}$
 - r^* $\{w \mid w \text{ is zero or more repetitions of a word in } r\}$

e.g.

$(a \mid b)b = \{ab, bb\}$

$a = \{a\}$ $a \mid b = \{a, b\}$
 $b = \{b\}$

$ab^* = \{a, ab, abb, abbb, abbbb, \dots\}$

$a^* = \{\epsilon, a, aa, aaa, aaaa, \dots\}$

$(aa \mid ab)^* = \{\epsilon, aa, ab, aaab, abaa, aaaa, abab, \dots\} \leftarrow aaaa aaaa ab$

$((aa)^* \mid (ab)^*) = \{aaaaaa, ababab, \dots\} \leftarrow aaaaaa, ababab, \dots$

Shorthands

$r? \equiv (r \mid \epsilon)$

$r^+ \equiv r(r^*)$

$\Sigma = \{a, b\}$

$ab^+a = \{aa, aba\}$

$1 + \dots r \mid 11 \quad 111 \quad 1111 \quad 11111 \quad \dots$

$$r^+ \equiv r(r^*)$$

$$ab^+a = \{aa, aba\}$$

$$b^+ = \{b, bb, bbb, bbbb, \dots\}$$

• REGULAR EXPRESSIONS & PROGRAMMING LANGUAGES

- Note: we cannot use regex to specify the "Language of all valid programs"

[What is the regex to validate HTML?]

closed

Why?

eg. $\{a^n b^n \mid n \in \mathbb{Z}\} = \{ab, aabb, aaabbb, \dots\}$

$$\{([])\} \quad \text{✓}$$

$$\{([])\} \quad \text{✗}$$

$$\langle a \rangle \dots \langle /a \rangle$$

$$\langle h3 \rangle \dots \langle /h3 \rangle$$

- Good enough to describe "tokens"

Token: a category of atomic strings.

lexeme: instance of a token.

E.G

while (flag == 62.3 else if) 14.7 != x + var_1

Keyword Ident oper lit real Keyword lit int oper Ident.

keywords = (while | if | else | (|) | do | return | ...)

Some tokens are infinite. \leftarrow integer literals
identifiers.

Σ = keyboard symbols (ASCII)

short-handly

$$[0-9] \equiv (0 | 1 | 2 | \dots | 9)$$

$$[A-Z] \equiv (A | B | C | D | \dots | Z)$$

$$[a-z] \equiv (a | b | c | d | \dots | z)$$

Regex for integer literals.

$(+|-)?([0-9])^+$

$\left\{ \begin{array}{l} -1, 123, -9, -00 \\ 000, +3, +000 \end{array} \right\}$

Regex for identifiers:

Fortran $([A-Z]|[a-z])([A-Z]|[a-z]|[0-9])^*$

Pascal $([A-Z]|[a-z])([A-Z]|[a-z]|[0-9]|' ')^* \{A, J, App13, 3d, \}$

Python $([A-Z]|[a-z]|' ')([A-Z]|[a-z]|[0-9]|' ')^* \{app-3, a-1-x \dots\}$

$\left\{ \begin{array}{l} \text{--init--}, \text{--name--}, \dots \\ \text{-----}, \end{array} \right\}$

we can write regex for:

- phone numbers
- urls
- serial no
- licence plates
- file names.

Regex for scientific notation

$1239e-5$

$(+|-)?[0-9]^+(\cdot[0-9]^+)? 'e' (+|-)?[0-9]^+$

$0.3e5$

— EOF —