

Regular Expressions



This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.

Outline and Objectives

Outline

- Motivation
- Regular expressions
 - Regular Ops.
 - Formal Definition
- Examples
- Applications

Learning Objectives

- Evaluate regular expressions.
- Construct regular expressions for simple languages.
- Describe regular operations.

Motivation

Task: Remove comments and white space from a text file. Assume comments start with two forward slashes (ie, //)

Python: `line = re.sub("//.*$", "", line)`

Perl: `$line =~ s//.*$//;`

Regular Expressions



Evaluating Expression

Arithmetic Expression

$$(5 + 3) * 2$$

- This evaluates to a number:
16
- Makes use of standard mathematical operations

Regular Expression

01*0

- This evaluates to a language: A string that starts and ends with '0' and has a continuous string of zero or more '1's
- Makes use of regular operations

Regular Operations

Regular operations operate on languages...

Let X and Y be regular languages, then following are regular operations. The result of each op. can be shown to be a regular language.

Concatenation: $X \bullet Y = \{ xy \mid x \in X \text{ AND } y \in Y \}$

Union: $X \cup Y = \{ x \mid x \in X \text{ OR } y \in Y \}$

Star: $X^* = \{ x_1 x_2 x_3 \dots x_i \mid i \geq 0 \text{ and } x_i \in X \}$

Regular Languages

R is a regular expression if R is any of the following...

- x for some x in an alphabet (set of valid chars.)
- ϵ or \emptyset ϵ is a language with empty string; \emptyset is an empty language
- $(R_1 \cup R_2)$, where R_1 and R_2 are regular languages
- $(R_1 \cdot R_2)$, where R_1 and R_2 are regular languages
- $(R_1)^*$, where R_1 is a regular language

Examples

Assume the valid set of characters is 0 & 1. Describe or explicitly state the language represented by each expression.

$(0 \cup 1)^*0 = \{ s \mid s \text{ ends in a } 0 \}$

$(1 \cup \epsilon)0 = \{ 10, 0 \}$

$(0 \cup 1)1(0 \cup 1) = \{ s \mid s \text{ contains 3 digits w/ the } 1 \text{ in the middle } \}$

Examples

Assume the valid set of characters is 0 & 1. Provide a regular expression that evaluates to the language described.

- A language that starts with 0 and ends with 1 and is 4 characters long. = $0(0 \cup 1)(0 \cup 1)1$

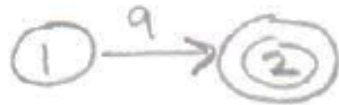
- A language that accepts any string containing a combination of 0's or 1's that has at least one character. = $(0 \cup 1)(0 \cup 1)^*$

- A string that has an even pair of 0's and no 1's. = $(00)^*$

Link to FSM

Regular expressions are equivalent to FSM. The idea is to construct a FSM that will accept any string from the language the RE represents.

a

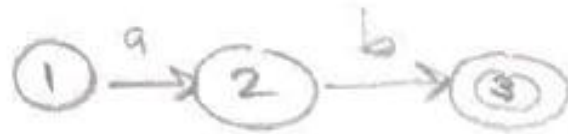


ϵ



Link to FSM

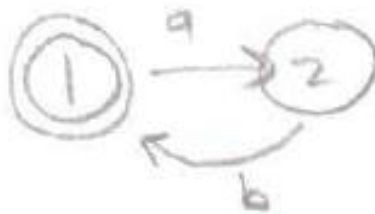
ab



a U b



(ab)*



Applications

You'll typically see regular expressions used for pattern matching.

When working with text strings, there are number of language specific extensions that make matching strings easier

Symbol	Meaning	Symbol	Meaning
^	Start of line	[]	Group of characters
\$	End of line	[^]	Group of characters to exclude
+	One or more	.	Any valid character

Examples

Provide a regular expression to match the following.

- White space at the beginning of a line
- The word 'add'
- Any number
- Any number with punctuation
- Whole line

`^[\t]*`

`add[\t\n]`

`[0-9]+`

`[0-9.,]+`

`^.*$`

You may need to escape the . to match a period instead of any character.

References

- ❑ Sipser M: *Introduction to the theory of computation*. 2nd ed. Boston: Thomson Course Technology; 2006.