

A crash course on regular expressions and regular grammars

Professor Hossein Saiedian

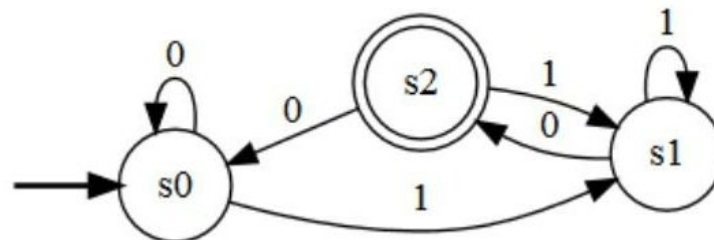
EECS 348: Software Engineering

Spring 2023

- A set of states, S
 - An initial state S_0
 - A possible set of one or more final states (F)
 - The initial and final states can be the same
- A set of inputs, I , also called the input vocabulary
- A set of outputs, O , also called the output vocabulary
- A finite-state automaton *accepts* a string x if $\delta(S_0, x) \in F$;
i.e., the finite state automaton ends up in a final state

FSM: Basic components

- A finite state automaton (FSA) is a graph with directed labeled arcs, two types of nodes, final and a unique start state
 - An FSA can have more than one final state
 - Also called a state machine
 - Deterministic FSA: For each state and for each input alphabet, there is exactly one transition
- The language accepted by a FSM is set of strings that move from start node to a final node



A couple of related terms

- **BNF** (Backus-Naur Form): a formal meta language used in defining context-free grammars, e.g., in compiler design
 - <terminal>
 - Defined as: $::=$
 - Alternative: $|$
 - Optional: $[]$
 - Combining: $()$
 - Maybe: $?$
 - Repeated: $(...)^*$ some also suggest $\{...\}$
 - One or more: $(...)^+$
 - Recursion: $\langle \text{integer} \rangle ::= \langle \text{digit} \rangle \mid \langle \text{integer} \rangle \langle \text{digit} \rangle$

Some symbols are just for convenience

Some BNF examples

NUMBER ::= WHOLE_PART FRACTION_PART?

WHOLE_PART ::= NUMERAL_SEQUENCE

FRACTION_PART ::= "." NUMERAL_SEQUENCE

NUMERAL_SEQUENCE ::= NUMERAL NUMERAL*

NUMERAL ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9

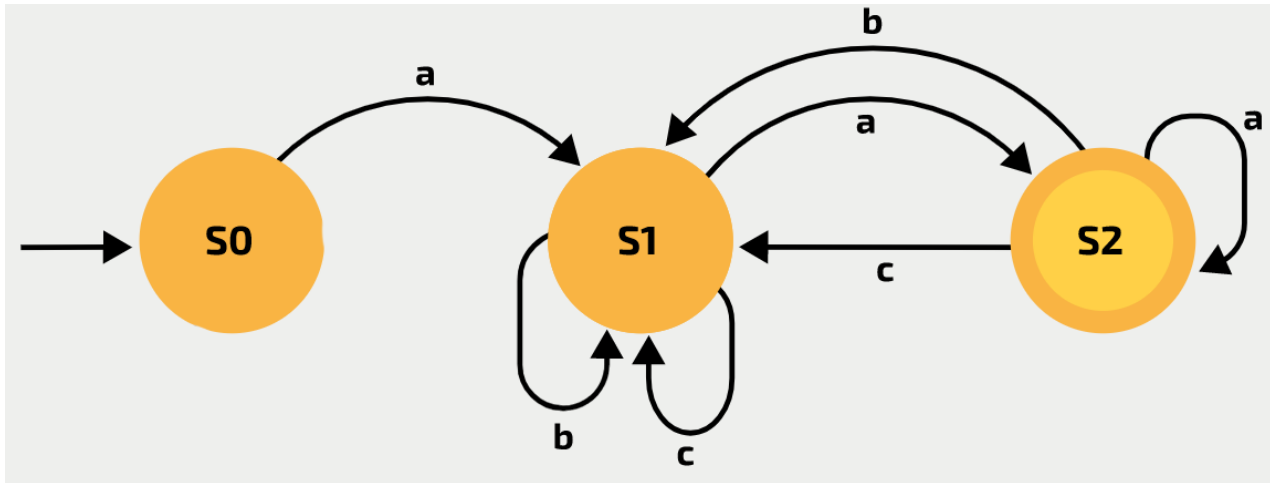
EXPRESSION ::= NUMERAL | (EXPRESSION OPTR EXPRESSION)

OPTR ::= + | -

- Not as powerful as BNF (e.g., no recursion)
- Regular expressions can express finite languages by defining a pattern for finite strings of symbols
- The grammar defined by regular expressions is known as regular grammar
- A regular expression is a pattern that can be recognized by a FSM
 - Regular expressions and FSMs are equivalent concepts
 - Semantically the same; one is textual (metal language), one is graphical

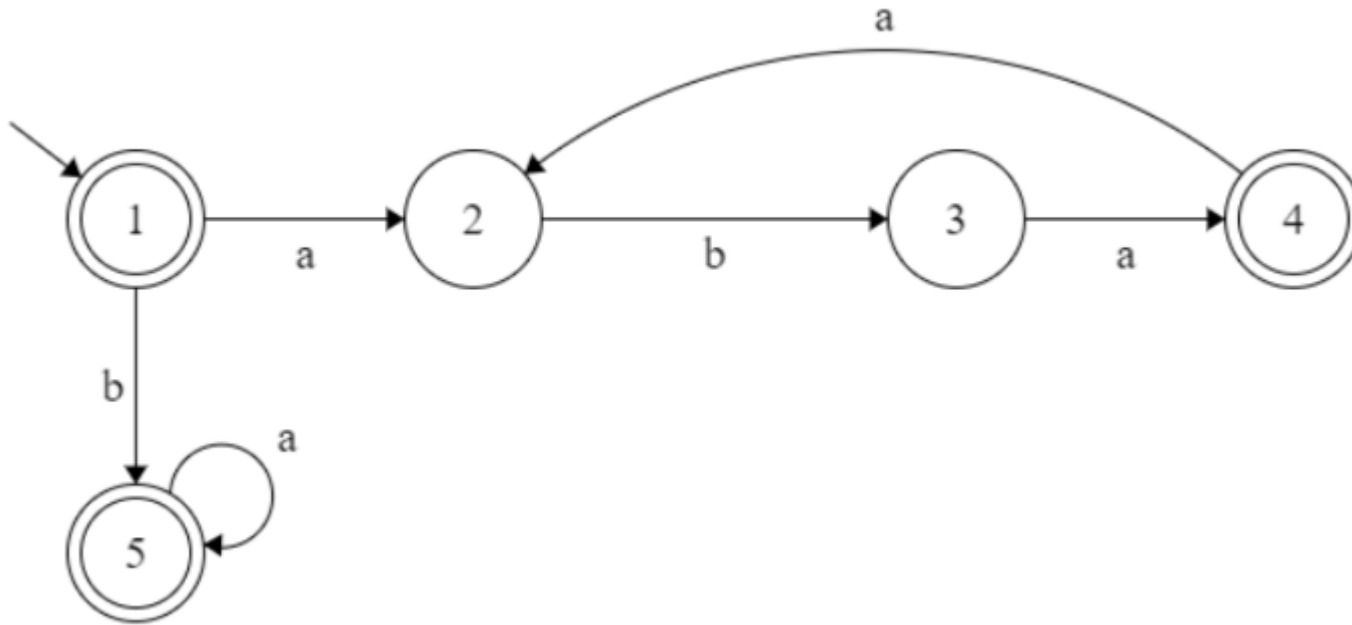
Examples of regular expressions

$a(a|b|c)^*a$



Examples of regular expressions

- $(aba)^+ | ba^*$



Examples of regular expressions

- What do these represent?

$L ::= \text{"A"} \mid \text{"B"} \mid \dots \mid \text{"a"} \mid \text{"b"} \mid \dots$

$D ::= 1 \mid 2 \mid \dots$

$L(\text{"_"} \mid L \mid D)^*$

$(\text{"+"} \mid \text{"-"})DD^*$

$(\text{"+"} \mid \text{"-"})D^+ \text{"."} D^+ [(\text{"E"} \mid \text{"e"})[\text{"+"} \mid \text{"-"}] D^+$

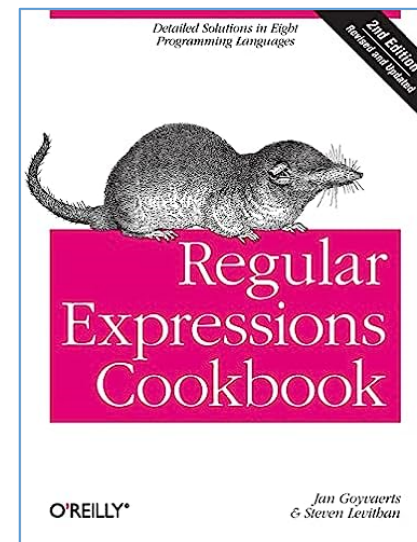
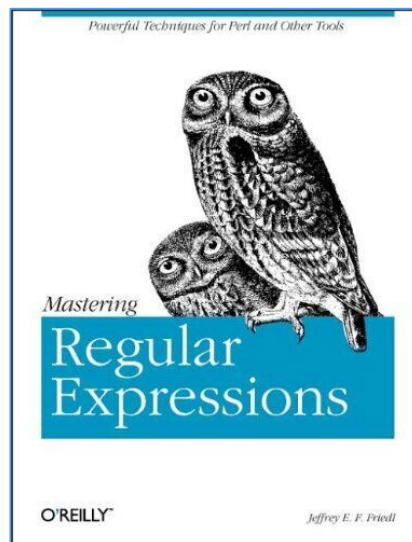
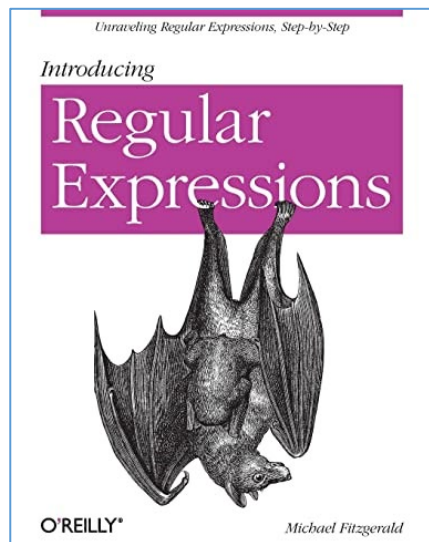
Regular expression in Unix

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)
- Common Unix symbols

Symbol	Descriptions
.	replaces any character
^	matches start of string
\$	matches end of string
*	matches up zero or more times the preceding character
\	Represent special characters
()	Groups regular expressions
?	Matches up exactly one character

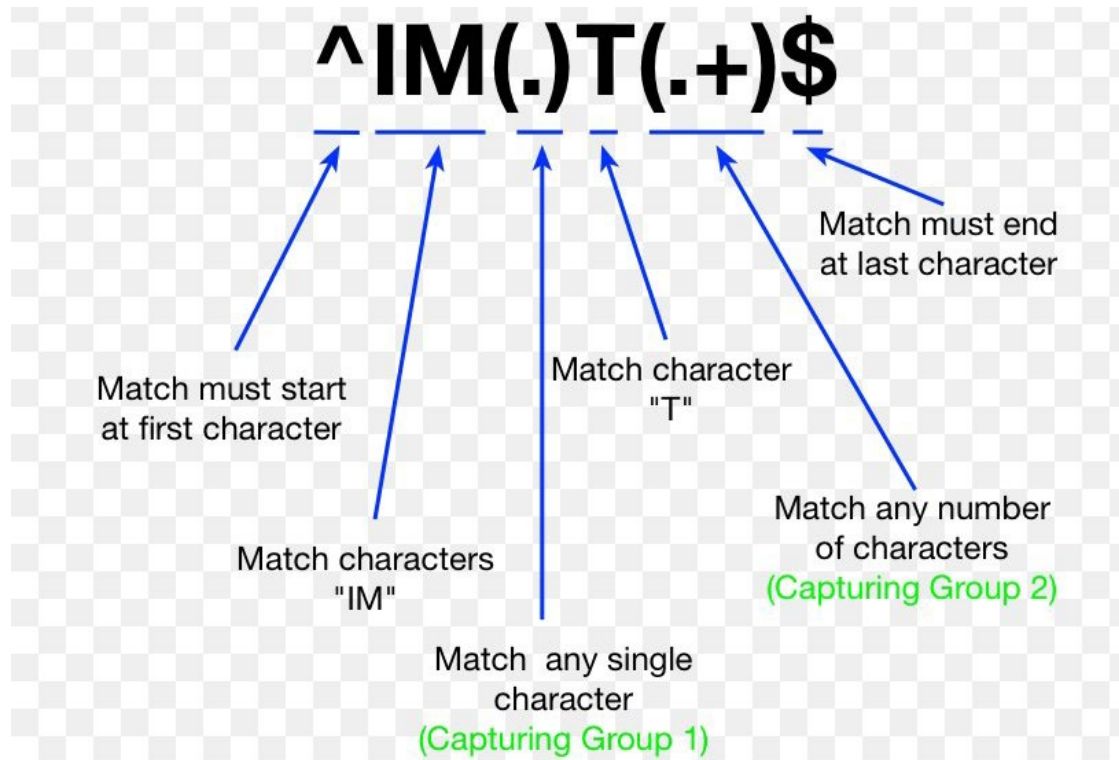
Regular expressions in Unix

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)



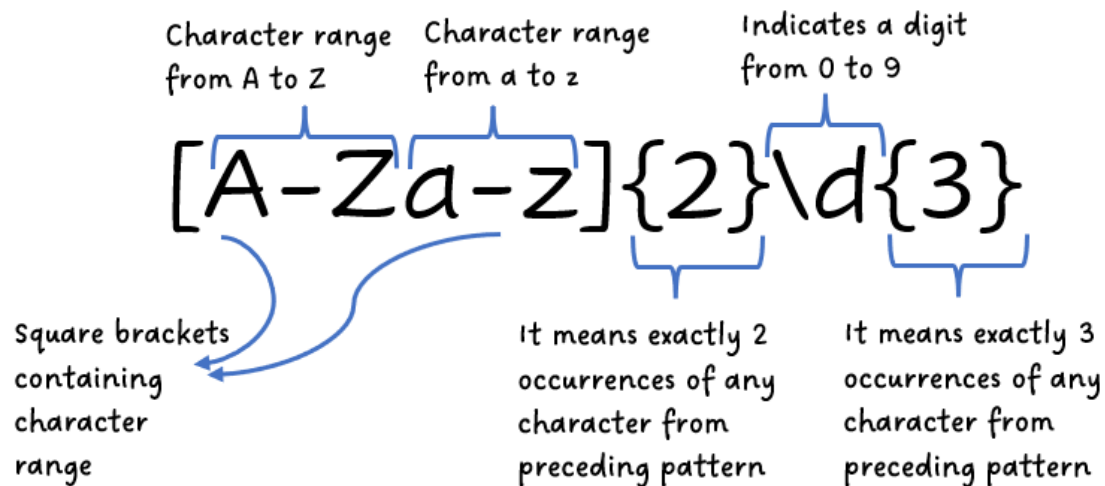
Regular expressions in Unix

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)



Regular expressions in Unix

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)

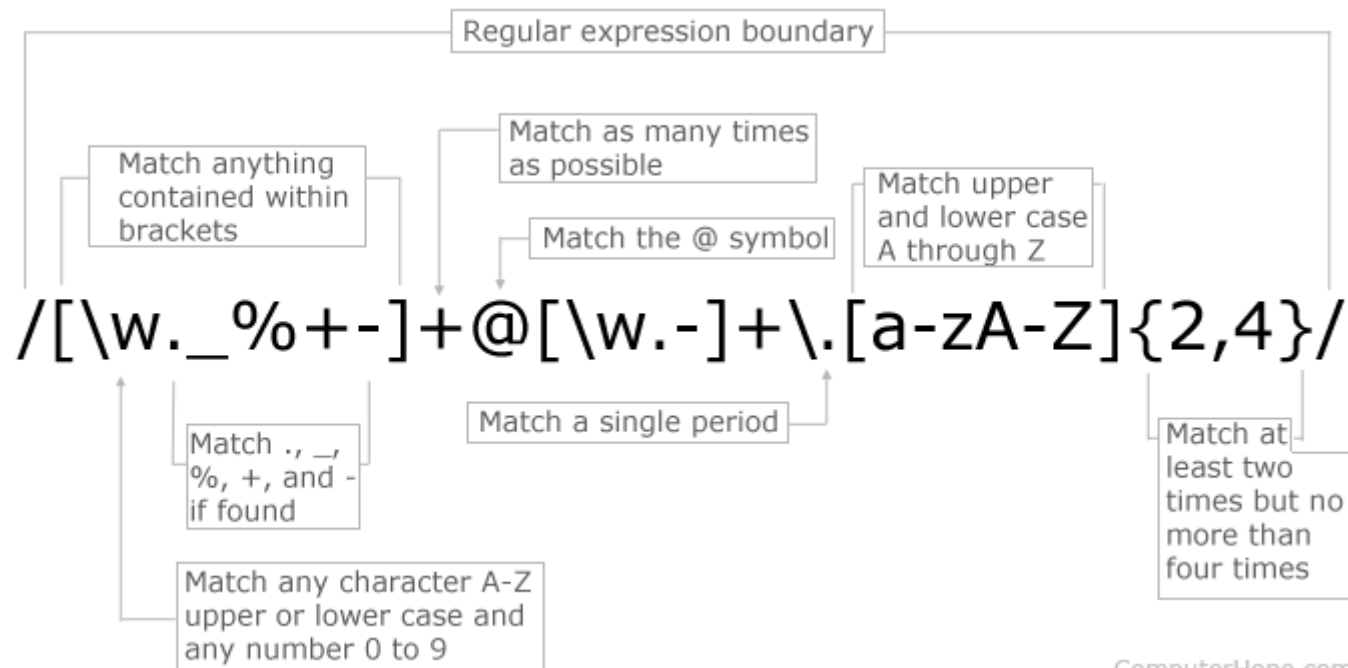


e.g., CS229, cs231

Examples that match above pattern

Regular expressions in Unix

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)

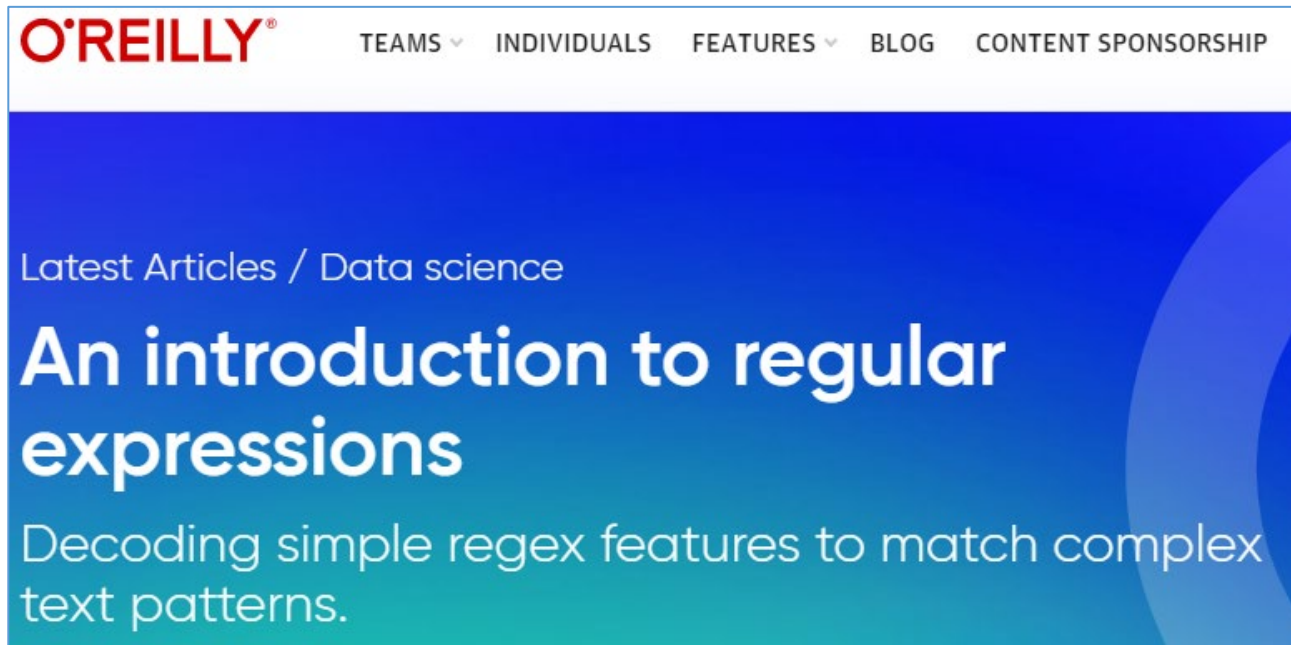


ComputerHope.com

A short RE cheat sheet

<code>[abc]</code>	A single character: a, b or c
<code>[^abc]</code>	Any single character but a, b, or c
<code>[a-z]</code>	Any single character in the range a-z
<code>[a-zA-Z]</code>	Any single character in the range a-z or A-Z
<code>^</code>	Start of line
<code>\$</code>	End of line
<code>\A</code>	Start of string
<code>\z</code>	End of string
<code>.</code>	Any single character
<code>\s</code>	Any whitespace character
<code>\S</code>	Any non-whitespace character
<code>\d</code>	Any digit
<code>\D</code>	Any non-digit
<code>\w</code>	Any word character (letter, number, underscore)
<code>\W</code>	Any non-word character
<code>\b</code>	Any word boundary character
<code>(...)</code>	Capture everything enclosed
<code>(a b)</code>	a or b
<code>a?</code>	Zero or one of a
<code>a*</code>	Zero or more of a
<code>a+</code>	One or more of a
<code>a{3}</code>	Exactly 3 of a
<code>a{3,}</code>	3 or more of a
<code>a{3,6}</code>	Between 3 and 6 of a

- A very similar concept, slightly different notation
- Extremely powerful for pattern matching (vim, grep, sed, ...)



<https://www.oreilly.com/content/an-introduction-to-regular-expressions>