
Regular Expressions

— Besant Technologies —

Regular Expression (Regex):

- A sequence of characters that defines a search pattern.
- Used for searching strings.
- Ex: `^a...s$` - any five letter strings starting with 'a' and ending with 's'

Expression	String	Matched?
<code>^a...s\$</code>	abs	No match
	alias	Match
	abyss	Match
	Alias	No match
	An abacus	No match

Metacharacters:

- Characters that are interpreted in a special way.
- List: [], ., ^, \$, *, +, ?, {}, (), \, |
-



- Specifies a set of characters to match.
- Ex: `[abc]` - looks for any of these characters to match.
- Can specify a range of characters using `-`
- Ex: `[a-z]`, `[0-9]` etc.,
- Can complement the character set using `^`
- Ex: `[^a-t]`, `[^1-5]` etc.,

Expression	String	Matched?
<code>[abc]</code>	a	1 match
	ac	2 matches
	Hey Jude	No match
	abc de ca	5 matches



- Matches any single character except new line.



Expression	String	Matched?
..	a	No match
	ac	1 match
	acd	1 match
	acde	2 matches (contains 4 characters)

^ - Caret

- Used to check if a string **starts** with certain character.



Expression	String	Matched?
^a	a	1 match
	abc	1 match
	bac	No match
^ab	abc	1 match
	acb	No match (starts with a but not followed by b)

\$ - dollar

- Used to check if any string ends with certain character/pattern.

Expression	String	Matched?
a\$	a	1 match
	formula	1 match
	cab	No match



- Matches zero or more occurrences of the pattern left to it.

Expression	String	Matched?
ma*n	mn	1 match
	man	1 match
	maaan	1 match
	main	No match (a is not followed by n)
	woman	1 match



- Matches one or more occurrences of the pattern left to it.



Expression	String	Matched?
ma+n	mn	No match (no a character)
	man	1 match
	maan	1 match
	main	No match (a is not followed by n)
	woman	1 match



- Matches zero or one occurrences of the pattern left to it.

Expression	String	Matched?
ma?n	mn	1 match
	man	1 match
	maaan	No match (more than one a character)
	main	No match (a is not followed by n)
	woman	1 match



- $\{n,m\}$ - means atleast n and atmost m repetitions of the pattern left to it.
-

Expression	String	Matched?
a{2,3}	abc dat	No match
	abc daat	1 match (at <u>daat</u>)
	aabc daaat	2 matches (at <u>aabc</u> and <u>daaat</u>)
	aabc daaaat	2 matches (at <u>aabc</u> and <u>daaaat</u>)

Expression	String	Matched?
[0-9]{2,4}	ab123csde	1 match (match at ab <u>123</u> csde)
	12 and 345673	2 matches (at <u>12</u> and <u>345673</u>)
	1 and 2	No match

| - or operator

Vertical bar `|` is used for alternation (`or` operator).

Expression	String	Matched?
<code>a b</code>	<code>cde</code>	No match
	<code>ade</code>	1 match (match at <code><u>a</u>de</code>)
	<code>acdbea</code>	3 matches (at <code><u>a</u><u>c</u><u>d</u><u>b</u><u>e</u><u>a</u></code>)

Here, `a|b` match any string that contains either `a` or `b`



- Used to group sub patterns.

Parentheses `()` is used to group sub-patterns. For example, `(a|b|c)xz` match any string that matches either `a` or `b` or `c` followed by `xz`

Expression	String	Matched?
<code>(a b c)xz</code>	<code>ab xz</code>	No match
	<code>abxz</code>	1 match (match at <code>ab<u>xz</u></code>)
	<code>axz cabxz</code>	2 matches (at <code><u>axz</u>bc cab<u>xz</u></code>)

\ - escape character

Backslash `\` is used to escape various characters including all metacharacters. For example,

`\$a` match if a string contains `$` followed by `a`. Here, `$` is not interpreted by a RegEx engine in a special way.

If you are unsure if a character has special meaning or not, you can put `\` in front of it. This makes sure the character is not treated in a special way.

Special Sequences:

- Most commonly used patterns.
- `'\A'`, `'\b'`, `'\B'`, `'\d'`, `'\D'`, `'\s'`, `'\S'`, `'\w'`, `'\W'`, `'\Z'`

\A:

`\A` - Matches if the specified characters are at the start of a string.

Expression	String	Matched?
<code>\Athe</code>	the sun	Match
	In the sun	No match

\Z:

`\Z` - Matches if the specified characters are at the end of a string.

Expression	String	Matched?
<code>\ZPython</code>	I like Python	1 match
	I like Python	No match
	Python is fun.	No match

\b:

`\b` - Matches if the specified characters are at the beginning or end of a word.

Expression	String	Matched?
\bfoo	football	Match
	a football	Match
	afootball	No match
foo\b	the foo	Match
	the afoo test	Match
	the afootest	No match

\B:

`\B` - Opposite of `\b`. Matches if the specified characters are **not** at the beginning or end of a word.

Expression	String	Matched?
\Bfoo	football	No match
	a football	No match
	afootball	Match
foo\b	the foo	No match
	the afoo test	No match
	the afootest	Match

\d:

`\d` - Matches any decimal digit. Equivalent to `[0-9]`

Expression	String	Matched?
<code>\d</code>	<code>12abc3</code>	3 matches (at <code><u>1</u>2abc<u>3</u></code>)
	<code>Python</code>	No match

\D:

`\D` - Matches any non-decimal digit. Equivalent to `[^0-9]`

Expression	String	Matched?
<code>\D</code>	1ab34"50	3 matches (at 1 <u>a</u> b34"50)
	1345	No match

\s:

`\s` - Matches where a string contains any whitespace character. Equivalent to `[\t\n\r\f\v]`.

Expression	String	Matched?
<code>\s</code>	Python RegEx	1 match
	PythonRegEx	No match

\S:

`\S` - Matches where a string contains any non-whitespace character. Equivalent to `[^ \t\n\r\f\v]`.

Expression	String	Matched?
<code>\S</code>	a b	2 matches (at <u>a</u> <u>b</u>)
		No match

\w:

`\w` - Matches any alphanumeric character (digits and alphabets). Equivalent to `[a-zA-Z0-9_]`. By the way, underscore `_` is also considered an alphanumeric character.

Expression	String	Matched?
<code>\w</code>	<code>12&" : ;c</code>	3 matches (at <code>12&" : ;c</code>)
	<code>%"> !</code>	No match

\W:

`\W` - Matches any non-alphanumeric character. Equivalent to `[^a-zA-Z0-9_]`

Expression	String	Matched?
<code>\W</code>	<code>1a2%c</code>	1 match (at <code>1<u>a</u><u>2</u><u>%</u><u>c</u></code>)
	<code>Python</code>	No match

re module:

- 'Re' module to work with regular expressions.
- <https://regex101.com/> - practice.
- Methods:
 - findall()
 - split()
 - sub()
 - subn()
 - match()
 - search()