



Regular Expressions and String Manipulation Playground

Sooner or later every programmer will use [Regular Expressions](#) in their code. To me they are hard to understand and reason... this post will try to clear up some of my confusion and provide examples of what these regular expressions look like.

Writing regular expressions

The first step is writing regular expressions. We have two ways of doing it: literals and constructors.

Regular expression literals

The simplest way to write regular expressions is to use them as a literal value in a constant or variable, the pattern is enclosed in slashes (/).

```
const myRE = /a+bc/;
```

Literal regular expressions are evaluated when the script is loaded. Use RegExp literals when you know what the regular expression will be.

Regular expression constructors

We can also build regular expressions using the RegExp constructor where we put the regular expression that we want to test in parenthesis.

```
const myRE = new RegExp('a+bc');
```

The constructors are evaluated when the script runs. Use constructors when you're not certain what the expression will be or when you're working with user input.

Modifying the regular expression: special characters and creating patterns

There are times when patterns like those we discussed are enough for our needs but there are times when we need to build more complicated expressions. We may want to match portions of a URL or make sure that the protocol used for the URL is `https://`.

There are special characters that we can use to enhance the expressions to perform specific tasks.

The table below, taken from MDN's [Writing a regular expression pattern](#) shows the characters you can use in regular expressions and what they do.

Special characters in regular expressions.

Character	Meaning
<code>\</code>	<p>Matches according to the following rules:</p> <ul style="list-style-type: none">▪ A backslash that precedes a non-special character indicates that the next character is special and is not to be interpreted literally. For example, a 'b' without a preceding '\' generally matches lowercase 'b's wherever they occur. But a '\b' by itself doesn't match any character; it denotes a word boundary.▪ A backslash that precedes a special character indicates that the next character is not special and should be interpreted literally. For example, the pattern <code>/a*/</code> relies on the special character '*' to match 0 or more a's. By contrast, the pattern <code>/a*/</code> denotes the '*' as not special, enabling matches with strings like 'a*'. <p>Do not forget to escape \ itself while using the <code>RegExp("pattern")</code> notation because \ is also an escape character in strings.</p>
<code>^</code>	Matches beginning of input. If the multiline flag is set to true, also matches immediately after a line break character.

Character	Meaning
	<p>For example, <code>/^A/</code> does not match the 'A' in "an A", but does match the 'A' in "An E".</p> <p>The '^' has a different meaning when it appears as the first character in a character set pattern.</p>
\$	<p>Matches end of input. If the multiline flag is set to true, also matches immediately before a line break character.</p> <p>For example, <code>/t\$/</code> does not match the 't' in "eater", but does match it in "eat".</p>
*	<p>Matches the preceding expression 0 or more times. Equivalent to <code>{0,}</code>.</p> <p>For example, <code>/bo*/</code> matches 'boooo' in "A ghost boooooed" and 'b' in "A bird warbled" but nothing in "A goat grunted".</p>
+	<p>Matches the preceding expression 1 or more times. Equivalent to <code>{1,}</code>.</p> <p>For example, <code>/a+/</code> matches the 'a' in "candy" and all the a's in "caaaaaaandy", but nothing in "cndy".</p>
?	<p>Matches the preceding expression 0 or 1 time. Equivalent to <code>{0, 1}</code>.</p> <p>For example, <code>/e?le?/</code> matches the 'el' in "angel" and the 'le' in "angle" and also the 'l' in "oslo".</p> <p>If used immediately after any of the quantifiers <code>*</code>, <code>+</code>, <code>?</code>, or <code>{}</code>, makes the quantifier non-greedy (matching the fewest possible characters), as opposed to the default, which is greedy (matching as many characters as possible). For example, applying <code>/\d+/</code> to "123abc" matches "123". But applying <code>/\d+?/</code> to that same string</p>

Character	Meaning
	<p>matches only the "1".</p> <p>Also used in lookahead assertions, as described in the <code>x(?=y)</code> and <code>x(?!y)</code> entries of this table.</p>
.	<p>(The decimal point) matches any single character except the newline character, by default.</p> <p>For example, <code>/ .n/</code> matches 'an' and 'on' in "nay, an apple is on the tree", but not 'nay'.</p> <p>If the <code>s ("dotAll")</code> flag is set to true, it also matches newline characters.</p>
(x)	<p>Matches 'x' and remembers the match, as the following example shows. The parentheses are called <i>capturing parentheses</i>.</p> <p>The <code>'(foo)'</code> and <code>'(bar)'</code> in the pattern <code>/(foo) (bar) \1 \2/</code> match and remember the first two words in the string "foo bar foo bar". The <code>\1</code> and <code>\2</code> denote the first and second parenthesized substring matches - foo and bar, matching the string's last two words. Note that <code>\1</code>, <code>\2</code>, ..., <code>\n</code> are used in the matching part of the regex. In the replacement part of a regex the syntax <code>\$1</code>, <code>\$2</code>, ..., <code>\$n</code> must be used, e.g.: <code>'bar foo'.replace(/(...) (...)/, '\$2 \$1')</code>. <code>&</code> means the whole matched string.</p>
(?:x)	<p>Matches 'x' but does not remember the match. The parentheses are called <i>non-capturing parentheses</i>, and let you define subexpressions for regular expression operators to work with. Consider the sample expression <code>/(?:foo){1,2}/</code>. If the expression was <code>/foo{1,2}/</code>, the <code>{1,2}</code> characters would apply only to the last 'o' in 'foo'. With the non-capturing parentheses, the <code>{1,2}</code> applies to the entire word 'foo'.</p>

Character	Meaning
<code>x(?=y)</code>	<p>Matches 'x' only if 'x' is followed by 'y'. This is called a lookahead.</p> <p>For example, <code>/Jack(?=Sprat)/</code> matches 'Jack' only if it is followed by 'Sprat'. <code>/Jack(?=Sprat Frost)/</code> matches 'Jack' only if it is followed by 'Sprat' or 'Frost'. However, neither 'Sprat' nor 'Frost' is part of the match results.</p>
<code>x(?!y)</code>	<p>Matches 'x' only if 'x' is not followed by 'y'. This is called a negated lookahead.</p> <p>For example, <code>/\d+(?!\.)/</code> matches a number only if it is not followed by a decimal point. The regular expression <code>/\d+(?!\.)/.exec("3.141")</code> matches '141' but not '3.141'.</p>
<code>(?<=y)x</code>	<p>Matches x only if x is preceded by y. This is called a lookbehind.</p> <p>For example, <code>/(?<=Jack)Sprat/</code> matches "Sprat" only if it is preceded by "Jack". <code>/(?<=Jack Tom)Sprat/</code> matches "Sprat" only if it is preceded by "Jack" or "Tom". However, neither "Jack" nor "Tom" is part of the match results.</p>
<code>(?<!=y)x</code>	<p>Matches x only if x is not preceded by y. This is called a negated lookbehind.</p> <p>For example, <code>/(?<!=)\d+/</code> matches a number only if it is not preceded by a minus sign. <code>/(?<!=)\d+/.exec('3')</code> matches "3". <code>/(?<!=)\d+/.exec('-3')</code> match is not found because the number is preceded by the minus sign.</p>

Character	Meaning
<code>x y</code>	<p>Matches 'x', or 'y' (if there is no match for 'x').</p> <p>For example, <code>/green red/</code> matches 'green' in "green apple" and 'red' in "red apple." The order of 'x' and 'y' matters. For example <code>a* b</code> matches the empty string in "b", but <code>b a*</code> matches "b" in the same string.</p>
<code>{n}</code>	<p>Matches exactly n occurrences of the preceding expression. N must be a positive integer.</p> <p>For example, <code>/a{2}/</code> doesn't match the 'a' in "candy," but it does match all of the a's in "caandy," and the first two a's in "caaandy."</p>
<code>{n,}</code>	<p>Matches at least n occurrences of the preceding expression. N must be a positive integer.</p> <p>For example, <code>/a{2,}/</code> will match "aa", "aaaa" and "aaaaa" but not "a"</p>
<code>{n,m}</code>	<p>Where n and m are positive integers and $n \leq m$. Matches at least n and at most m occurrences of the preceding expression. When m is omitted, it's treated as ∞.</p> <p>For example, <code>/a{1,3}/</code> matches nothing in "cndy", the 'a' in "candy," the first two a's in "caandy," and the first three a's in "caaaaaaandy". Notice that when matching "caaaaaaandy", the match is "aaa", even though the original string had more a's in it.</p>
<code>[xyz]</code>	<p>Character set. This pattern type matches any one of the characters in the brackets, including escape sequences. Special characters like the</p>

Character	Meaning
	<p>dot(.) and asterisk(*) are not special inside a character set, so they don't need to be escaped. You can specify a range of characters by using a hyphen, as the following examples illustrate.</p> <p>The pattern [a-d], which performs the same match as [abcd], matches the 'b' in "brisket" and the 'c' in "city". The patterns /[a-z.]+/ and /[\\w.]+/ match the entire string "test.i.ng".</p>
[^xyz]	<p>A negated or complemented character set. That is, it matches anything that is not enclosed in the brackets. You can specify a range of characters by using a hyphen. Everything that works in the normal character set also works here.</p> <p>For example, [^abc] is the same as [^a-c]. They initially match 'r' in "brisket" and 'h' in "chop."</p>
[\\b]	<p>Matches a backspace (U+0008). You need to use square brackets if you want to match a literal backspace character. (Not to be confused with \\b.)</p>
\\b	<p>Matches a <i>word boundary</i>. A word boundary matches the position between a word character followed by a non-word character, or between a non-word character followed by a word character, or the beginning of the string, or the end of the string. A word boundary is not a "character" to be matched; like an anchor, a word boundary is not included in the match. In other words, the length of a matched word boundary is zero. (Not to be confused with [\\b].)</p> <p>Examples using the input string "moon": /\\bm/ matches, because the '\\b' is at the beginning of the string;</p> <p>the '\\b' in /oo\\b/ does not match, because the '\\b' is both preceded and followed by word characters;</p>

Character	Meaning
	<p>the '\b' in /oon\b/ matches, because it appears at the end of the string;</p> <p>the '\b' in /\w\b\w/ will never match anything, because it is both preceded and followed by a word character..</p> <p>Note: JavaScript's regular expression engine defines a specific set of characters to be "word" characters. Any character not in that set is considered a non-word character. This set of characters is fairly limited: it consists solely of the Roman alphabet in both upper- and lower-case, decimal digits, and the underscore character. Accented characters, such as "é" or "ü" are, unfortunately, treated as non-word characters for the purposes of word boundaries, as are ideographic characters in general.</p>
\B	<p>Matches a non-<i>word boundary</i>. This matches the following cases:</p> <ul style="list-style-type: none"> ▪ Before the first character of the string. ▪ After the last character of the string,. ▪ Between two word characters ▪ Between two non-word characters ▪ The empty string <p>For example, /\B. ./ matches 'oo' in "noonday", and /y\B. / matches 'ye' in "possibly yesterday."</p>
\cX	<p>Where <i>X</i> is a character ranging from A to Z. Matches a control character in a string.</p> <p>For example, /\cM/ matches control-M (U+000D) in a string.</p>
\d	Matches a digit character. Equivalent to [0-9].

Character	Meaning
	For example, <code>/\d/</code> or <code>/[0-9]/</code> matches '2' in "B2 is the suite number."
<code>\D</code>	Matches a non-digit character. Equivalent to <code>[^0-9]</code> . For example, <code>/\D/</code> or <code>/[^0-9]/</code> matches 'B' in "B2 is the suite number."
<code>\f</code>	Matches a form feed (U+000C).
<code>\n</code>	Matches a line feed (U+000A).
<code>\r</code>	Matches a carriage return (U+000D).
<code>\s</code>	Matches a white space character, including space, tab, form feed, line feed. Equivalent to <code>[\f\n\r\t\v\u00a0\u1680\u2000-\u200a\u2028\u2029\u202f\u205f\u3000\ufe0f]</code> . For example, <code>/\s\w*/</code> matches ' bar' in "foo bar."
<code>\S</code>	Matches a character other than white space. Equivalent to <code>[^ \f\n\r\t\v\u00a0\u1680\u2000-\u200a\u2028\u2029\u202f\u205f\u3000\ufe0f]</code> . For example, <code>/\S*/</code> matches 'foo' in "foo bar."
<code>\t</code>	Matches a tab (U+0009).
<code>\v</code>	Matches a vertical tab (U+000B).
<code>\w</code>	Matches any alphanumeric character including the underscore. Equivalent to <code>[A-Za-z0-9_]</code> .

Character	Meaning
	For example, <code>/\w/</code> matches 'a' in "apple," '5' in "\$5.28," and '3' in "3D."
<code>\W</code>	Matches any non-word character. Equivalent to <code>[^A-Za-z0-9_]</code> . For example, <code>/\W/</code> or <code>/[^A-Za-z0-9_]/</code> matches '%' in "50%."
<code>\n</code>	Where <i>n</i> is a positive integer, a back reference to the last substring matching the <i>n</i> parenthetical in the regular expression (counting left parentheses). For example, <code>/apple(,)\sorange\1/</code> matches 'apple, orange,' in "apple, orange, cherry, peach."
<code>\0</code>	Matches a NULL (U+0000) character. Do not follow this with another digit, because <code>\0<digits></code> is an octal (base 8) sequence. Instead use <code>\x00</code> .
<code>\xhh</code>	Matches the character with the code hh (two hexadecimal digits)

Going back to finding a regular expression that tests if our URL is to a secure site we can do the following:

```
const secureURL = /https:\/\/;
```

We'll use the regular expression when we look at how to test using regular expressions in the next section.

Additional flags for regular expressions

One last set of flags to worry about. These are additional parameters for the expression that will allow additional flexibility beyond what special characters allow.

The table below show what these additional flags are and what they do.

Regular expression flags	
Flag	Description
g	Global search.
i	Case-insensitive search.
m	Multi-line search.
s	Allows . to match newline characters.
u	"Treat a pattern as a sequence of unicode code points
y	Perform a "sticky" search that matches starting at the current position in the target string.

The flags are appended to the end of the regular expression and are a part of it (no adding or removing flags after the regular expression) and they will change the way we build the regular expressions we want to use.

The literal expression now looks like this:

```
const re = /pattern/flags;
```

And the constructor changes to the one below

```
const re = new RegExp('pattern', 'flags');
```

Matching text against regular expressions

We have multiple ways to test strings against our regular expressions. Which one to use depends on what results you want to accomplish and how much information about the matches and the results you want to keep and use.

The methods are listed below:

Methods that use regular expressions

Method	Description
<code>exec</code>	A <code>RegExp</code> method that executes a search for a match in a string. It returns an array of information or null on a mismatch.
<code>test</code>	A <code>RegExp</code> method that tests for a match in a string. It returns true or false.
<code>match</code>	A <code>String</code> method that executes a search for a match in a string. It returns an array of information or null on a mismatch.
<code>matchAll</code>	A <code>String</code> method that returns an iterator containing all of the matches, including capturing groups.
<code>search</code>	A <code>String</code> method that tests for a match in a string. It returns the index of the match, or -1 if the search fails.
<code>replace</code>	A <code>String</code> method that executes a search for a match in a string, and replaces the matched substring with a replacement substring.
<code>split</code>	A <code>String</code> method that uses a regular expression or a fixed string to break a string into an array of substrings.

What happens when the string has Unicode Characters?

Unicode special characters for regular expressions.

Character	Meaning
<code>\uhhhh</code>	Matches the character with the code hhhh (four hexadecimal digits).
<code>\u{hhhh}</code>	(only when u flag is set) Matches the character with the Unicode value hhhh (hexadecimal digits).

Making life easier: groups

Matching multiple items in the string

[String.prototype.matchAll](#)

A common use case of global (g) or sticky (y) regular expressions is applying it

to a string and iterating through all of the matches. The new `String.prototype.matchAll` API makes this easier than ever before, especially for regular expressions with capture groups:

```
const string = 'Favorite GitHub repos: tc39/ecma262 v8/v8.dev';
const regex = /\b(?:<owner>[a-z0-9]+\)\b(?:<repo>[a-z0-9\.\.]+\b/g;

for (const /\b(?:<owner>[a-z0-9]+\)\b(?:<repo>[a-z0-9\.\.]+\b/g(`${match[0]}`) {
  console.log(`→ owner: ${match.groups.owner}`);
  console.log(`→ repo: ${match.groups.repo}`);
}

`→ owner: ${match.groups.owner}` // Output:
//
// tc39/ecma262 at 23 with 'Favorite GitHub repos: tc39/ecma262 v8/v8.dev'
// → owner: tc39
// → repo: ecma262
// v8/v8.dev at 36 with 'Favorite GitHub repos: tc39/ecma262 v8/v8.dev'
// → owner: v8
// → repo: v8.dev
```

Links, resources and Ideas for further experimentation

- <https://alligator.io/js/regular-expressions-for-regular-people/>
- <https://www.smashingmagazine.com/2019/02/regexp-features-regular-expressions/>