

Strings and Regular Expressions

Relevant Links

- Flanagan's book, sections 3.2, 10.1, 10.2, (10.3 optional)
- MDN's guide on RegExps¹
- MDN's RegExp reference²
- MDN's String object reference³

Strings in Javascript

- String literals are formed by surrounding the text in single or double quotes: "a string", 'another string'.
- In ES 6, there is a third string literal, called *template literal*, and we will discuss those further later.
- Escape characters are combinations that hold special meaning. Refer to table 3.1 in the book, or in MDN's String reference.
- The plus operator is overloaded to cause string concatenation. Other values will be coerced to strings if necessary.

```
"hi" + 0;      // "hi0"  
0 + "hi";      // "0hi"  
"1" + 0;       // "10"
```

Note: Oftentimes you read numeric input from a text field. You must convert it into a number first, to avoid examples like the above.

- Strings can be accessed via array indexing: "hi there"[4] == "h", "hello".length == 5. There is no separate type for single characters.
- You can iterate over the characters in a string using for example a for-of loop.
- Strings are immutable: You cannot change their value, you can only create a new string.

String methods

charAt⁴ An alternative for accessing a specific location in the string

charCodeAt⁵ Returns the numeric Unicode value of the character at a specific location.

concat⁶ Concatenates strings

¹https://developer.mozilla.org/en-US/docs/Web/JavaScript/Guide/Regular_Expressions

²https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/RegExp

³https://developer.mozilla.org/en-US/docs/Web/JavaScript/Reference/Global_Objects/String

indexOf⁷ Search within the string for the first occurrence of a substring.
lastIndexOf⁸ Similar to indexOf, but finds the last occurrence instead.
split⁹ Split a string into an array of strings based on a separator parameter.
substr¹⁰ Extract a section of a string
substring¹¹ Very similar to substr, but slightly different arguments
toLowerCase¹² Converts to lower case
toUpperCase¹³ Converts to upper case
startsWith¹⁴ Checks if the string starts with the provided string
endsWith¹⁵ Checks if the string ends with the provided string
includes¹⁶ Checks if the string includes with the provided string
repeat¹⁷ Repeats the string a number of times

Regular Expressions in Javascript

- Regular expressions are patterns that allow us to express very intelligent searches in strings.
- The full language of regular expressions is quite rich and will take some time to fully digest. We will only scratch the surface here.
- A regular expression literal is marked by two forward slashes: `/stuffHere/`. The second slash may optionally be followed by the letters “i”, “g”, and/or “m”, which affect how the expression behaves.
 - i** the regular expression ignores case, so an “a” would match both lower case and upper case A’s.
 - g** the regular expression becomes “global”, so it will keep on searching for more matches after the first one.
 - m** the regular expression possibly matches across multiple lines.
 - u** offers some Unicode-related features.
- Between the two forward slashes, a number of elements can be present, each with its meaning. Here are the most important ones:
 - c** Any character except for a few special characters that we will mention further down matches itself. This includes numbers, letters, some punctuation etc.
 - \t** Matches a tab character.
 - \n** Matches a newline character.
 - \r** Matches a carriage return. Often appears before the newline character to denote end of line.
 - .** A dot matches any character at all, except for newlines. Precede the dot with a backslash if you want to match a literal dot.
 - \d** Matches a digit.
 - \D** Matches any non-digit.
 - \w** Matches an alphanumeric character or underscore.
 - \W** Matches any character that is not alphanumeric or underscore.
 - \s** Matches a whitespace character (spaces, tabs, newlines and a couple more).

\S Matches a non-whitespace character.

**** Matches an actual backslash.

\c For a character that has a special meaning by itself, escape it with a backslash if you want its literal meaning. e.g. `\\.`, `\\?`, `*`, `\\|`.

[xyz] Matches any of the characters inside the brackets. You can also indicate ranges, like `[a-z.0-9]`, which matches any lowercase character or digit or dot.

[^xyz] Matches anything BUT what follows the caret.

^ Doesn't match an actual character, but it must match the beginning of a line.

\$ Doesn't match an actual character, but it must match the end of a line.

\b Doesn't match an actual character, but it must match a word boundary. Useful for matching whole words. For instance `/\\bcat\\b/` will match in "12cat" but not in "cats".

(...) Treats the contents as a group. In many regular expression searches, you would be able to access the matched pieces. Also called *capturing parentheses*.

\\n where n is a number. Matches a previously matched group. For example `(.)\\1` will match an occurrence of the same character back-to-back.

(?:...) *Non-capturing parentheses*. Useful when you simply want to group some elements together (to impose some operator precedence for example), but do not care for capturing the result separately.

x* Matches the expression x zero or more times.

x+ Matches the expression one or more times.

x? Matches the expression zero or one time.

x*? Matches zero or more times, but in a "non-greedy" fashion, meaning it will try to match the smallest possible number of times.

x|y Matches either x or y.

- Some examples would be in order:

```
/([\\w\\-]+\\d+)@hanover\\.edu/
```

This will match all student emails at Hanover (faculty/staff emails don't end in digits). It looks for a word character or dash, one or more times, followed by one or more digits, then a literal at sign, "hanover", then a literal dot, then "edu".

```
/()(.)().\\3\\2\\1/
```

This regular expression will match palindromes with total length 6, e.g. something like "abccba". The parenthesized dots match one arbitrary element each, then the backticked numbers refer to those matches in reverse order.

```
/((.)(.)(.))\\4\\3\\2\\1/
```

QUESTION: What does the above regex match?

```
/^\\s*[+-]?(?:\\d+\\.?[\\d\\.\\d+)?\\d*(?:[eE][+-]?\\d+)?\\s*$/
```

This is a fairly complicated regular expression that matches numbers in decimal or scientific format. It starts with a caret, and ends in a dollar sign, meaning that the contents must match the entire line/string. It allows for an arbitrary number of whitespace characters on either end. After that, it has an optional sign, followed by: either at least one digit followed by an optional dot, or a dot following

```
/"(?:\\|\"|[\^"])*"/
```

Using regular expressions

- Use the matches of `r` as separators to split the string at.

- When working with regular expressions, it is always a good idea to create a list of test strings first, then see how the regular expression behaves on those strings.