

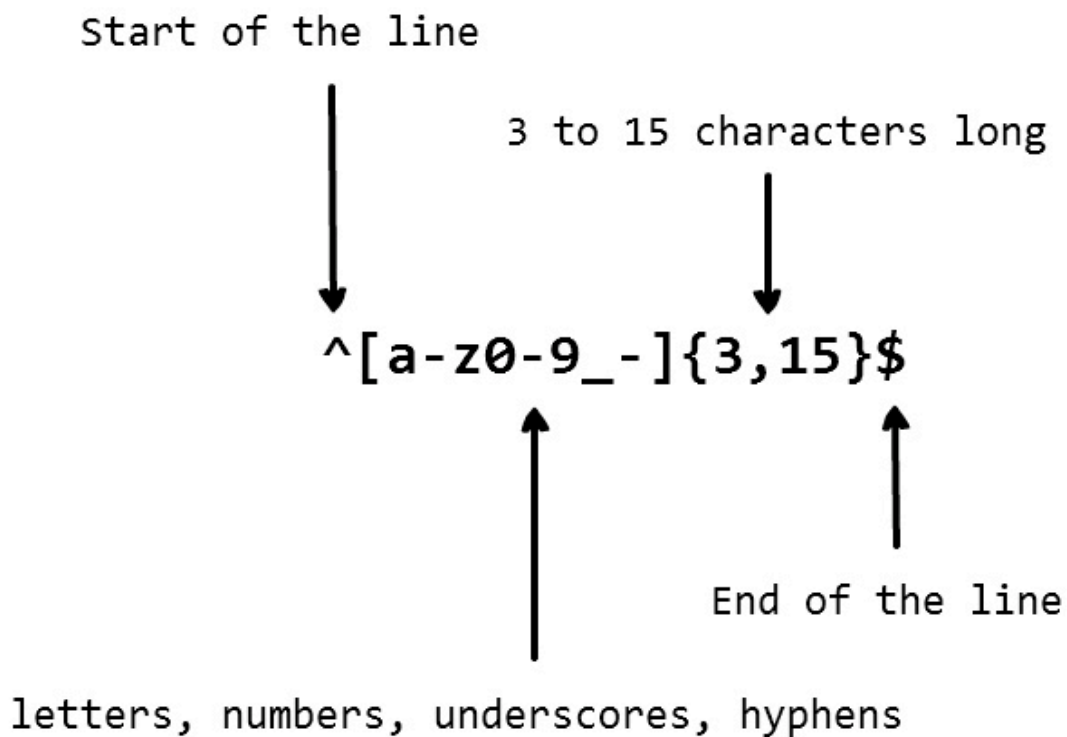
Regular Expression

Regular expression is a group of characters or symbols which is used to find a specific pattern from some text.

Regular expression is used for replacing a text within a string, validating form, extract a substring from a string based upon a pattern match, and so much more.

Eg:

We want to allow the username to contain letters, numbers, underscores and hyphens. We also want to limit the number of characters in username so it does not look ugly.



Regular expression

Meta character	Description
.	Period matches any single character except a line break.
[]	Character class. Matches any character contained between the square brackets.
[^]	Negated character class. Matches any character that is not contained between the square brackets
*	Matches 0 or more repetitions of the preceding symbol.
+	Matches 1 or more repetitions of the preceding symbol.
?	Makes the preceding symbol optional.
{n,m}	Braces. Matches at least "n" but not more than "m" repetitions of the preceding symbol.
(xyz)	Character group. Matches the characters xyz in that exact order.
	Alternation. Matches either the characters before or the characters after the symbol.
\	Escapes the next character. This allows you to match reserved characters [] () { } . * + ? ^ \$ \
^	Matches the beginning of the input.
\$	Matches the end of the input.

2.1 Full stop

Full stop . is the simplest example of meta character. The meta character . matches any single character. It will not match return or newline characters.

For example, the regular expression .ar means: any character, followed by the letter a, followed by the letter r.

.ar => The **car** parked in the **garage**.

2.2 Character set

Character sets are also called character class. Square brackets are used to specify character sets. Use a hyphen inside a character set to specify the characters' range. The order of the character range inside square brackets

doesn't matter. For example, the regular expression [Tt]he means: an uppercase T or lowercase t, followed by the letter h, followed by the letter e.

[Tt]he => **The** car parked in **the** garage.

A period inside a character set, however, means a literal period. The regular expression `ar[.]` means: a lowercase character `a`, followed by letter `r`, followed by a period `.` character.

`ar[.]` => A garage is a good place to park a **car**.

2.2.1 Negated character set

In general, the caret symbol represents the start of the string, but when it is typed after the opening square bracket it negates the character set. For example, the regular expression `[^c]ar` means: any character except `c`, followed by the character `a`, followed by the letter `r`.

```
[^c]ar => The car parked in the garage.
```

2.3 Repetitions

Following meta characters `+`, `*` or `?` are used to specify how many times a subpattern can occur. These meta characters act differently in different situations.

2.3.1 The Star

The symbol `*` matches zero or more repetitions of the preceding matcher. The regular expression `a*` means: zero or more repetitions of preceding lowercase character `a`. But if it appears after a character set or class then it finds the repetitions of the whole character set. For example, the regular expression `[a-z]*` means: any number of lowercase letters in a row.

```
[a-z]* => The car parked in the garage #21.
```

The `*` symbol can be used with the meta character `.` to match any string of characters `.*`. The `*` symbol can be used with the whitespace character `\s` to match a string of whitespace characters. For example, the expression `\sca\s` means: zero or more spaces, followed by lowercase character `c`, followed by lowercase character `a`, followed by lowercase character `t`, followed by zero or more spaces.

```
\s*cat\s* => The fat cat sat on the concatenation.
```

2.3.2 The Plus

The symbol `+` matches one or more repetitions of the preceding character. For example, the regular expression `c.+t` means: lowercase letter `c`, followed by at least one character, followed by the lowercase character `t`. It needs to be clarified that `t` is the last `t` in the sentence.

```
c.+t => The fat cat sat on the mat.
```

2.3.3 The Question Mark

In regular expression the meta character `?` makes the preceding character optional. This symbol matches zero or one instance of the preceding character. For example, the regular expression `[T]?he` means: Optional the uppercase letter `T`, followed by the lowercase character `h`, followed by the lowercase character `e`.

```
[T]he => The car is parked in the garage.
```

2.4 Braces

In regular expression braces that are also called quantifiers are used to specify the number of times that a character or a group of characters can be repeated. For example, the regular expression `[0-9]{2,3}` means: Match at least 2 digits but not more than 3 (characters in the range of 0 to 9).

```
[0-9]{2,3} => The number was 9.9997 but we rounded it off to 10.0.
```

We can leave out the second number. For example, the regular expression `[0-9]{2,}` means: Match 2 or more digits. If we also remove the comma the regular expression `[0-9]{3}` means: Match exactly 3 digits.

```
[0-9]{2,} => The number was 9.9997 but we rounded it off to 10.0.
```

```
[0-9]{3} => The number was 9.9997 but we rounded it off to 10.0.
```

2.5 Character Group

Character group is a group of sub-patterns that is written inside Parentheses (...). As we discussed before that in regular expression if we put a quantifier after a character then it will repeat the preceding character. But if we put quantifier after a character group then it repeats the whole character group. For example, the regular expression (ab)* matches zero or more repetitions of the character ab. We can also use the alternation | meta character inside character group. For example, the regular expression (c|g|p)ar means: lowercase character c, g or p, followed by character a, followed by character r.

```
(c|g|p)ar => The car is parked in the garage.
```

2.6 Alternation

In regular expression Vertical bar | is used to define alternation. Alternation is like a condition between multiple expressions. Now, you may be thinking that character set and alternation works the same way. But the big difference between character set and alternation is that character set works on character level but alternation works on expression level. For example, the regular expression (T|t)he|car means: uppercase character T or lowercase t, followed by lowercase character h, followed by lowercase character e or lowercase character c, followed by lowercase character a, followed by lowercase character r.

```
(T|t)he|car => The car is parked in the garage.
```

2.7 Escaping special character

Backslash \ is used in regular expression to escape the next character. This allows us to specify a symbol as a matching character including reserved characters { } [] / \ + * . \$ ^ | ?. To use a special character as a matching character prepend \ before it.

For example, the regular expression . is used to match any character except newline. Now to match . in an input string the regular expression (f|c|m)at.? means: lowercase

letter f, c or m followed by lowercase

character a, followed by lowercase letter t, followed by optional . character.

```
(f|c|m)at\.? => The fat cat sat on the mat.
```

2.8 Anchors

In regular expressions, we use anchors to check if the matching symbol is the starting symbol or ending symbol of the input string. Anchors are of two types: First type is Caret ^ that check if the matching character is the start character of the input and the second type is Dollar \$ that checks if matching character is the last character of the input string.

2.8.1 Caret

Caret ^ symbol is used to check if matching character is the first character of the input string. If we apply the following regular expression ^a (if a is the starting symbol) to input string abc it matches a. But if we apply regular expression ^b on above input string it does not match anything. Because in input string abc b is not the starting symbol. Let's take a look at another regular expression ^{T|t}he which means: uppercase character T or lowercase character t is the start symbol of the input string, followed by lowercase character h, followed by lowercase character e.

```
(T|t)he => The car is parked in the garage.
```

```
^(T|t)he => The car is parked in the garage.
```

2.8.2 Dollar

Dollar

symbol is used to check if matching character is the last character of the input string. For example, means: a lowercase character a, followed by lowercase character t, followed by a . character and the matcher must be end of the string.

```
(at\.) => The fat cat. sat. on the mat.
```

3. Shorthand Character Sets

Regular expression provides shorthands for the commonly used character sets, which offer convenient shorthands for commonly used regular expressions. The shorthand character sets are as follows:

Shorthand	Description
.	Any character except new line
\w	Matches alphanumeric characters: <code>[a-zA-Z0-9_]</code>
\W	Matches non-alphanumeric characters: <code>[^\w]</code>
\d	Matches digit: <code>[0-9]</code>
\D	Matches non-digit: <code>[^\d]</code>
\s	Matches whitespace character: <code>[\t\n\f\r\p{Z}]</code>
\S	Matches non-whitespace character: <code>[^\s]</code>

4. Lookaround

Lookbehind and lookahead (also called lookaround) are specific types of non-capturing groups (Used to match the pattern but not included in matching list). Lookaheads are used when we have the condition that this pattern is preceded or followed by another certain pattern. For example, we want to get all numbers that are preceded by

character from the following input string 4.44 and

10.88. We will use following regular expression `(? <= $)[0-9\.] *` which means : get all the character. Following are the lookarounds that are used in regular expressions:

Symbol	Description
?=	Positive Lookahead
?!	Negative Lookahead
?<=	Positive Lookbehind
?<!	Negative Lookbehind

4.1 Positive Lookahead

The positive lookahead asserts that the first part of the expression must be followed by the lookahead expression. The returned match only contains the text that is matched by the first part of the expression. To define a positive lookahead, parentheses are used. Within those parentheses, a question mark with equal sign is used like this: `(?=...)`. Lookahead expression is written after the equal sign inside parentheses. For example, the regular expression `(T|t)he(=?\sfat)` means: optionally match lowercase letter t or uppercase letter T, followed by letter h, followed by letter e. In parentheses we define positive lookahead which tells regular expression engine to match The or the which are followed by the word fat.

```
(T|t)he(=?\sfat) => The fat cat sat on the mat.
```

4.2 Negative Lookahead

Negative lookahead is used when we need to get all matches from input string that are not followed by a pattern. Negative lookahead is defined same as we define positive

lookahead but the only difference is instead of equal = character we use negation ! character i.e. (?!...). Let's take a look at the following regular expression (T|t)he(?!\sfat) which means: get all The or the words from input string that are not followed by the word fat precedes by a space character.

```
(T|t)he(?!\sfat) => The fat cat sat on the mat.
```

4.3 Positive Lookbehind

Positive lookbehind is used to get all the matches that are preceded by a specific pattern. Positive lookbehind is denoted by (?<=...). For example, the regular expression (?<=(T|t)he\s)(fat|mat) means: get all fat or mat words from input string that are after the word The or the.

```
(?<=(T|t)he\s)(fat|mat) => The fat cat sat on the mat.
```

4.4 Negative Lookbehind

Negative lookbehind is used to get all the matches that are not preceded by a specific pattern. Negative lookbehind is denoted by (?<!). For example, the regular expression (?<!(T|t)he\s)(cat) means: get all cat words from input string that are not after the word The or the.

```
(?<!(T|t)he\s)(cat) => The cat sat on cat.
```

5. Flags

Flags are also called modifiers because they modify the output of a regular expression. These flags can be used in any order or combination, and are an integral part of the RegExp.

Flag	Description
i	Case insensitive: Sets matching to be case-insensitive.
g	Global Search: Search for a pattern throughout the input string.
m	Multiline: Anchor meta character works on each line.

5.1 Case Insensitive

The `i` modifier is used to perform case-insensitive matching. For example, the regular expression `/The/gi` means: uppercase letter `T`, followed by lowercase character `h`, followed by character `e`. And at the end of regular expression the `i` flag tells the regular expression engine to ignore the case. As you can see we also provided `g` flag because we want to search for the pattern in the whole input string.

```
The => The fat cat sat on the mat.
```

```
/The/gi => The fat cat sat on the mat.
```

5.2 Global search

The `g` modifier is used to perform a global match (find all matches rather than stopping after the first match). For example, the regular expression `/(at)/g` means: any character except new line, followed by lowercase character `a`, followed by lowercase character `t`. Because we provided `g` flag at the end of the regular expression now it will find all matches in the input string, not just the first one (which is the default behavior).

```
/(at)/ => The fat cat sat on the mat.
```

```
/.(at)/g => The fat cat sat on the mat.
```

5.3 Multiline

The m modifier is used to perform a multi-line match. As we discussed earlier

anchors (

)are used to check if pattern is the beginning of the input or end of the input string. But if we u

/gm means: lowercase

character a, followed by lowercase character t, optionally anything except

new line. And because of m flag now regular expression engine matches pattern at the end of

each line in a string.

```
/.at(.)?$/ => The fat  
             cat sat  
             on the mat.
```

```
"/.at(.)?$/gm" => The fat  
                 cat sat  
                 on the mat.
```

6. Greedy vs lazy matching

By default regex will do greedy matching , means it will match as long as

possible. We can use ? to match in lazy way means as short as possible

```
/(.*at)/ => The fat cat sat on the mat.
```

```
/(.*?at)/ => The fat cat sat on the mat.
```

Digits:

```
/^[0-9]+$/
```

Alphabetic Characters

```
/^[a-zA-Z]+$/
```

Alpha-Numeric Characters

```
/^[a-zA-Z0-9]+$/
```

Date (MM/DD/YYYY)/(MM-DD-YYYY)/(MM.DD.YYYY)/(MM DD YYYY)

```
/^(0?[1-9]|1[012])[- /.](0?[1-9]|[12][0-9]|3[01])[- /.](19|20)?[0-9]{2}$/
```

Email Address

```
/^[a-zA-Z0-9._%~]+@[a-zA-Z0-9.-]+.[a-zA-Z]{2,4}$/
```

Password

The password must contain one lowercase letter, one uppercase letter, one number, one unique character such as !@#\$%^&? and be at least 6 characters long.

```
/^(?=.*{6,})(?=.*d)(?=.*[A-Z])(?=.*[a-z])(?=.*[!@#$%^&*? ]).*$/
```

US Phone Numbers

```
/^(?([0-9]{3}))?[-. ]?( [0-9]{3})[-. ]?( [0-9]{4})$/
```

US Zip code

```
/^[0-9]{5}(:-[0-9]{4})?$/
```

URL

```
/(((( [A-Za-z]{3,9}:(?:\\\/\\\/)?(?:[-;:&=\\+\$,\\w]+@)?[A-Za-z0-9.-]+|(? :www. | [-;:&=\\+\$,\\w]+@) [A-Za-z0-9.-]+) (((?:\\\/[\\+~%\\\/.\\w_-]*)?)??(?:[-\\+=&;%@.\\w_]*)#?(?:[\\w]*)?)?) /
```

IP Address

```
/^(?:(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?).){3}(?:25[0-5]|2[0-4][0-9]|[01]?[0-9][0-9]?)$/
```

HTML Tag

```
/^<([a-z]+)([<]+)*(?:>(.*)</1>|s+>)$/
```

