

Perl Compatible Regular Expressions in a Nutshell

All features on this cheat sheet are available in latest Perl, however, the implementation and UTF-8 support in other environments may differ slightly, e.g. Ruby supports named groups and UTF-8 only since version 1.9.1.

Syntax

Regular expressions are usually wrapped in forward slashes (e.g. `/match me/`), however, you can use almost any other character for instance to prevent heavy quoting when matching against file paths (e.g. `%/bin/grep%` is the same as `/\bin\grep/`).

Perl:

```
$string =~ m/pattern/modifier
$string =~ s/pattern/replace/flags
```

<http://www.perl.com/doc/manual/html/pod/perlre.html>

Ruby:

```
string.match(/pattern/flags) or string.match(/pattern/flags) {|match| ... }
string.sub(/pattern/flags, replace) or string.sub(/pattern/flags) {|replace| ... }
string.gsub(/pattern/flags, replace) or string.gsub(/pattern/flags) {|replace| ... }
```

<http://www.apidock.com/ruby/String>

Python:

```
import re
re.search("(?flags)pattern", string)
re.sub("(?flags)pattern", replace, string, count)
```

<http://docs.python.org/library/re.html>

PHP:

```
preg_match('/pattern/flags', $string)
preg_replace('/pattern/flags', 'replace', $string)
```

<http://www.php.net/manual/en/ref.pcre.php>

Java:

```
string.matches("(?flags)pattern")
string.replaceFirst("(?flags)pattern", "replace")
string.replaceAll("(?flags)pattern", "replace")
```

<http://java.sun.com/javase/6/docs/api/java/util/regex/package-summary.html>

JavaScript:

```
string.match(/pattern/flags)
string.replace(/pattern/flags, "replace")
```

https://developer.mozilla.org/en/Core_JavaScript_1.5_Guide/Regular_Expressions

C++: (requires Boost.Regex library)

```
boost::regex_search(string, boost::regex("pattern", "flags"))
boost::regex_replace(string, boost::regex("pattern", "flags"), "replace")
```

http://www.boost.org/doc/libs/1_42_0/libs/regex/doc/html

Objective-C: (requires RegexKit or RegexKit Lite)

```
[string isMatchedByRegex:@"(?flags)pattern"]
[string stringByMatching:@"(?flags)pattern" replace:1 withString:@"replace"]
[string stringByMatching:@"(?flags)pattern" replace:RKReplaceAll withString:@"replace"]
```

<http://regexkit.sourceforge.net/> <http://regexkit.sourceforge.net/RegexKitLite>

Shell:

```
grep -P "(?flags)pattern" file.txt
```

<http://www.gnu.org/software/grep/>

Characters

These are the usual suspects well known from any C-ish language:

a	match the character a
3	match the number 3
\$a or #{a}	match the contents of a variable \$a (e.g. Perl) or a (e.g. Ruby) respectively
\n	newline (NL, LF)
\r	return (CR)
\f	form feed (FF)
\t	tab (TAB)
\x3C	character with the hex code 3C
\u561A	character with the hex code 561A
\e	escape character (alias \u001B)
\c...	control character

Wildcards

Wildcards match if a character belongs to the designated class of characters:

.	match any character (use \. to match a single dot, → quoting)
\w	alphanumeric + underscore (shortcut for [0-9a-zA-Z_])
\W	any character not covered by \w
\d	numeric (shortcut for [0-9])
\D	any character not covered by \d
\s	whitespace (shortcut for [\t\n\r\f])
\S	any character not covered by \s
[...]	any character listed: [a5!*d-g] means the characters a, 5, !, * and d, e, f, g
[^...]	any character not listed: [^a5!*d-g] means anything but the characters a, 5, !, * and d, e, f, g

Boundaries

Boundaries match the spots between characters and therefore have no width of their own (also called zero-width, → extensions):

\b	matches at a word boundary (spot between \w and \W)
\B	matches anything but a word boundary

- ^** matches at the beginning of a line (m) or entire string (s)
- \A** matches at the beginning of the entire string **!! RUBY: USE THIS INSTEAD OF ^ !!**
- \$** matches at the end of a line (m) or entire string (s)
- \Z** matches at the end of the entire string ignoring a trailing \n
- \z** matches at the end of the entire string **!! RUBY: USE THIS INSTEAD OF \$!!**
- \G** matches where the previous regex call left off (→ flag g)

Grouping

Any of the above constructs can be grouped to improve readability and create a reference for use in **pattern** or **replace** (→ replacing):

(...) the group is assigned to the references **\1** and **\1** as well as **\$1** (outside of the regex context)

(...) (...) etc first group is **\1**, **\1** and **\$1**, second group is **\2**, **\2** and **\$2** etc

(...|...|...) matches if one of the group options matches and assigns it to **\1**, **\1** and **\$1**

(?<name>...) the group is named **name** and assigned to the references **\g<name>** and **\k<name>** (Python uses **(P?<name>...)** instead)

You can use named groups to make complex regular expressions much more readable. The following example will match IPv4 entries from a hosts file, e.g. "myhost 192.168.1.1". Whitespaces are ignored due to the **x** flag (→ flags), the **{0}** on the first three lines in effect turns those named groups into mere placeholders and the actual pattern is but on the fourth line:

```
/(?<host> [a-z.-]+ ) {0}
 (?<byte> \d{1,3} ) {0}
 (?<ip> ( (?<byte>\. ) {3} \g<byte> ) {0}
 \g<ip>\s+\g<host>/x
```

Named group matches are available outside of the regex context as well, however, this is of course implemented differently for each programming language. Here is an example for Ruby 1.9.1:

```
m = "aabb33dd".match(/(?<numbers>\d+)/)
puts m[:numbers] # => "33"
```

Extensions

Less common functionality is covered by extensions using the **(?...) syntax**. Extensions do **not** create a reference like grouping does.

- (?:...|...)** same as grouping, but no reference is created
- (?=...)** zero-width positive lookahead assertion
- (?!...)** zero-width negative lookahead assertion
- (?<=...)** zero-width positive lookbehind assertion (no quantifiers allowed within)
- (?<!...)** zero-width negative lookbehind assertion (no quantifiers allowed within)
- (?>...)** zero-width independent subexpression
- (? (...) | ...)** conditional expression
- (?flags)** apply the flag(s) within the current group from this point forward (→ flags)
- (?flags:...)** apply the flag(s) for this pattern (no backreference created!)
- (?#...)** zero-width comment (no round brackets allowed in comment text)

Quantifiers

Most of the above constructs may be quantified by adding one of the following symbols after them:

- ?** match 1 or 0 times
- *** 0 or more times
- +** 1 or more times
- {n}** exactly *n* times
- {n,}** at least *n* times
- {n,m}** at least *n* but not more than *m* times, as often as possible

Greediness

This is a very important feature, ignore it and you are destined to produce clumsy and error prone regex! Quantifiers are greedy by default which means they match as often as possible. Limit their hunger by adding a **?** after them. Here's an example applied to the title of this section:

- G.*e** matches as "Greedyne" (**.*** halts on the *last* subsequent **e**)
- G.*?e** matches as "Gre" (**.*** halts on the *first* subsequent **e**)

Quoting

Characters which carry functional meaning in a given context are called meta characters and have to be quoted in order to match against the literal character. Unnecessary quoting is less likely to cause trouble than not quoting where it is required.

- \...** quote a single meta character: ***** matches a star instead of acting as a quantifier
- \Q ... \E** ignore all meta characters in between

Replacing

The following symbols have special meanings in the **replace** part:

- \1, \2 etc** include the contents of the corresponding group (→ grouping)
- \{1}000** same as the above, use curly brackets if numbers follow the symbol
- \l** lowercase the following character
- \L ... \E** lowercase all characters in between
- \u** uppercase the following character
- \U ... \E** uppercase all characters in between

Flags

Optional flags determine the behaviour of the regex as a whole. May be used within the **(?flags)** construct (→ extensions):

- i** case-insensitive pattern matching
- m** multiple lines: **.** does not match **\n** (Ruby uses this per default)
- s** single line: **.** matches **\n** (Ruby uses **m** for this instead)
- x** ignore whitespaces in pattern for better readability

The following **cannot** be used within the **(?flags)** construct:

- g** apply the regex as many times as possible (i.e. for global replace)
- e** evaluate the **replace** part as if it were source code **!! DANGEROUS !!**
- o** compile the pattern only once and therefore perform variable substitutions only once