

Pattern Syntax

(No version information available, might be only in CVS)

Pattern Syntax — Describes PCRE regex syntax

Description

The PCRE library is a set of functions that implement regular expression pattern matching using the same syntax and semantics as Perl 5, with just a few differences (see below). The current implementation corresponds to Perl 5.005.

Differences From Perl

The differences described here are with respect to Perl 5.005.

1. By default, a whitespace character is any character that the C library function `isspace()` recognizes, though it is possible to compile PCRE with alternative character type tables. Normally `isspace()` matches space, formfeed, newline, carriage return, horizontal tab, and vertical tab. Perl 5 no longer includes vertical tab in its set of whitespace characters. The `\v` escape that was in the Perl documentation for a long time was never in fact recognized. However, the character itself was treated as whitespace at least up to 5.002. In 5.004 and 5.005 it does not match `\s`.
2. PCRE does not allow repeat quantifiers on lookahead assertions. Perl permits them, but they do not mean what you might think. For example, `(?!a){3}` does not assert that the next three characters are not "a". It just asserts that the next character is not "a" three times.
3. Capturing subpatterns that occur inside negative lookahead assertions are counted, but their entries in the offsets vector are never set. Perl sets its numerical variables from any such patterns that are matched before the assertion fails to match something (thereby succeeding), but only if the negative lookahead assertion contains just one branch.
4. Though binary zero characters are supported in the subject string, they are not allowed in a pattern string because it is passed as a normal C string, terminated by zero. The escape sequence `"\x00"` can be used in the pattern to represent a binary zero.
5. The following Perl escape sequences are not supported: `\l`, `\u`, `\L`, `\U`. In fact these are implemented by Perl's general string-handling and are not part of its pattern matching engine.
6. The Perl `\G` assertion is not supported as it is not relevant to single pattern matches.
7. Fairly obviously, PCRE does not support the `(?{code})` and `(??{code})` construction. However, there is support for recursive patterns.
8. There are at the time of writing some oddities in Perl 5.005_02 concerned with the settings of captured strings when part of a pattern is repeated. For example, matching "aba" against the pattern `^(a(b)?)+$/` sets `$2` to the value "b", but matching "aabbaa" against `^(aa(bb)?)+$/` leaves `$2` unset. However, if the pattern is changed to `^(aa(b(b)?)+$/` then `$2` (and `$3`) get set. In Perl 5.004 `$2` is set in both cases, and that is also **TRUE** of PCRE. If in the future Perl changes to a consistent state that is different, PCRE may change to follow.
9. Another as yet unresolved discrepancy is that in Perl 5.005_02 the pattern `^(a)?(1a|b)+$/` matches the string "a", whereas in PCRE it does not. However, in both Perl and PCRE `^(a)?a/`

matched against "a" leaves \$1 unset.

10. PCRE provides some extensions to the Perl regular expression facilities:

- a. Although lookbehind assertions must match fixed length strings, each alternative branch of a lookbehind assertion can match a different length of string. Perl 5.005 requires them all to have the same length.
- b. If [PCRE_DOLLAR_ENDONLY](#) is set and [PCRE_MULTILINE](#) is not set, the \$ meta-character matches only at the very end of the string.
- c. If [PCRE_EXTRA](#) is set, a backslash followed by a letter with no special meaning is faulted.
- d. If [PCRE_UNGREEDY](#) is set, the greediness of the repetition quantifiers is inverted, that is, by default they are not greedy, but if followed by a question mark they are.

Regular Expression Details

Introduction

The syntax and semantics of the regular expressions supported by PCRE are described below. Regular expressions are also described in the Perl documentation and in a number of other books, some of which have copious examples. Jeffrey Friedl's "Mastering Regular Expressions", published by O'Reilly (ISBN 1-56592-257-3), covers them in great detail. The description here is intended as reference documentation.

A regular expression is a pattern that is matched against a subject string from left to right. Most characters stand for themselves in a pattern, and match the corresponding characters in the subject. As a trivial example, the pattern *The quick brown fox* matches a portion of a subject string that is identical to itself.

Meta-characters

The power of regular expressions comes from the ability to include alternatives and repetitions in the pattern. These are encoded in the pattern by the use of *meta-characters*, which do not stand for themselves but instead are interpreted in some special way.

There are two different sets of meta-characters: those that are recognized anywhere in the pattern except within square brackets, and those that are recognized in square brackets. Outside square brackets, the meta-characters are as follows:

| | |
|----|--|
| \ | general escape character with several uses |
| ^ | assert start of subject (or line, in multiline mode) |
| \$ | assert end of subject (or line, in multiline mode) |
| . | match any character except newline (by default) |
| [| start character class definition |
|] | end character class definition |
| / | |

| | |
|---|---|
| (| start of alternative branch |
| (| start subpattern |
|) | end subpattern |
| ? | extends the meaning of (, also 0 or 1 quantifier, also quantifier minimizer |
| * | 0 or more quantifier |
| + | 1 or more quantifier |
| { | start min/max quantifier |
| } | end min/max quantifier |

Part of a pattern that is in square brackets is called a "character class". In a character class the only meta-characters are:

| | |
|---|---|
| \ | general escape character |
| ^ | negate the class, but only if the first character |
| - | indicates character range |
| / | terminates the character class |

The following sections describe the use of each of the meta-characters.

Backslash

The backslash character has several uses. Firstly, if it is followed by a non-alphanumeric character, it takes away any special meaning that character may have. This use of backslash as an escape character applies both inside and outside character classes.

For example, if you want to match a "*" character, you write "*" in the pattern. This applies whether or not the following character would otherwise be interpreted as a meta-character, so it is always safe to precede a non-alphanumeric with "\" to specify that it stands for itself. In particular, if you want to match a backslash, you write "\\".

Note: Single and double quoted PHP [strings](#) have special meaning of backslash. Thus if \ has to be matched with a regular expression \, then "\\\" or '\\\\' must be used in PHP code.

If a pattern is compiled with the [PCRE_EXTENDED](#) option, whitespace in the pattern (other than in a character class) and characters between a "#" outside a character class and the next newline character are ignored. An escaping backslash can be used to include a whitespace or "#" character as part of the pattern.

A second use of backslash provides a way of encoding non-printing characters in patterns in a visible

manner. There is no restriction on the appearance of non-printing characters, apart from the binary zero that terminates a pattern, but when a pattern is being prepared by text editing, it is usually easier to use one of the following escape sequences than the binary character it represents:

| | |
|-------------------|---|
| <code>\a</code> | alarm, that is, the BEL character (hex 07) |
| <code>\cx</code> | "control-x", where x is any character |
| <code>\e</code> | escape (hex 1B) |
| <code>\f</code> | formfeed (hex 0C) |
| <code>\n</code> | newline (hex 0A) |
| <code>\r</code> | carriage return (hex 0D) |
| <code>\t</code> | tab (hex 09) |
| <code>\xhh</code> | character with hex code hh |
| <code>\ddd</code> | character with octal code ddd, or backreference |

The precise effect of `"\cx"` is as follows: if "x" is a lower case letter, it is converted to upper case. Then bit 6 of the character (hex 40) is inverted. Thus `"\cz"` becomes hex 1A, but `"\c{"` becomes hex 3B, while `"\c;"` becomes hex 7B.

After `"\x"`, up to two hexadecimal digits are read (letters can be in upper or lower case). In *UTF-8 mode*, `"\x{...}"` is allowed, where the contents of the braces is a string of hexadecimal digits. It is interpreted as a UTF-8 character whose code number is the given hexadecimal number. The original hexadecimal escape sequence, `\xhh`, matches a two-byte UTF-8 character if the value is greater than 127.

After `"\0"` up to two further octal digits are read. In both cases, if there are fewer than two digits, just those that are present are used. Thus the sequence `"\0\x\07"` specifies two binary zeros followed by a BEL character. Make sure you supply two digits after the initial zero if the character that follows is itself an octal digit.

The handling of a backslash followed by a digit other than 0 is complicated. Outside a character class, PCRE reads it and any following digits as a decimal number. If the number is less than 10, or if there have been at least that many previous capturing left parentheses in the expression, the entire sequence is taken as a *back reference*. A description of how this works is given later, following the discussion of parenthesized subpatterns.

Inside a character class, or if the decimal number is greater than 9 and there have not been that many capturing subpatterns, PCRE re-reads up to three octal digits following the backslash, and generates a single byte from the least significant 8 bits of the value. Any subsequent digits stand for themselves. For example:

| | |
|-------------------|-----------------------------------|
| <code>\040</code> | is another way of writing a space |
| <code>\40</code> | |

is the same, provided there are fewer than 40 previous capturing subpatterns

`\7`

is always a back reference

`\11`

might be a back reference, or another way of writing a tab

`\011`

is always a tab

`\0113`

is a tab followed by the character "3"

`\113`

is the character with octal code 113 (since there can be no more than 99 back references)

`\377`

is a byte consisting entirely of 1 bits

`\81`

is either a back reference, or a binary zero followed by the two characters "8" and "1"

Note that octal values of 100 or greater must not be introduced by a leading zero, because no more than three octal digits are ever read.

All the sequences that define a single byte value can be used both inside and outside character classes. In addition, inside a character class, the sequence `"\b"` is interpreted as the backspace character (hex 08). Outside a character class it has a different meaning (see below).

The third use of backslash is for specifying generic character types:

`\d`

any decimal digit

`\D`

any character that is not a decimal digit

`\h`

any horizontal whitespace character (since PHP 5.2.4)

`\H`

any character that is not a horizontal whitespace character (since PHP 5.2.4)

`\s`

any whitespace character

`\S`

any character that is not a whitespace character

`\v`

any vertical whitespace character (since PHP 5.2.4)

`\V`

any character that is not a vertical whitespace character (since PHP 5.2.4)

`\w`

any "word" character

`\W`

any "non-word" character

Each pair of escape sequences partitions the complete set of characters into two disjoint sets. Any given character matches one, and only one, of each pair.

A "word" character is any letter or digit or the underscore character, that is, any character which can be part of a Perl "word". The definition of letters and digits is controlled by PCRE's character tables, and

may vary if locale-specific matching is taking place. For example, in the "fr" (French) locale, some character codes greater than 128 are used for accented letters, and these are matched by `\w`.

These character type sequences can appear both inside and outside character classes. They each match one character of the appropriate type. If the current matching point is at the end of the subject string, all of them fail, since there is no character to match.

The fourth use of backslash is for certain simple assertions. An assertion specifies a condition that has to be met at a particular point in a match, without consuming any characters from the subject string. The use of subpatterns for more complicated assertions is described below. The backslashed assertions are

| | |
|-----------------|--|
| <code>\b</code> | word boundary |
| <code>\B</code> | not a word boundary |
| <code>\A</code> | start of subject (independent of multiline mode) |
| <code>\Z</code> | end of subject or newline at end (independent of multiline mode) |
| <code>\z</code> | end of subject (independent of multiline mode) |
| <code>\G</code> | first matching position in subject |

These assertions may not appear in character classes (but note that `"\b"` has a different meaning, namely the backspace character, inside a character class).

A word boundary is a position in the subject string where the current character and the previous character do not both match `\w` or `\W` (i.e. one matches `\w` and the other matches `\W`), or the start or end of the string if the first or last character matches `\w`, respectively.

The `\A`, `\Z`, and `\z` assertions differ from the traditional circumflex and dollar (described below) in that they only ever match at the very start and end of the subject string, whatever options are set. They are not affected by the [PCRE_MULTILINE](#) or [PCRE_DOLLAR_ENDONLY](#) options. The difference between `\Z` and `\z` is that `\Z` matches before a newline that is the last character of the string as well as at the end of the string, whereas `\z` matches only at the end.

The `\G` assertion is true only when the current matching position is at the start point of the match, as specified by the `offset` argument of [preg_match\(\)](#). It differs from `\A` when the value of `offset` is non-zero. It is available since PHP 4.3.3.

`\Q` and `\E` can be used to ignore regexp metacharacters in the pattern since PHP 4.3.3. For example: `\w+\Q.\$ \E$` will match one or more word characters, followed by literals `.$` and anchored at the end of the string.

`\K` can be used to reset the match start since PHP 5.2.4. For example, the pattern `foo\Kbar` matches "foobar", but reports that it has matched "bar". The use of `\K` does not interfere with the setting of captured substrings. For example, when the pattern `(foo)\Kbar` matches "foobar", the first substring is still set to "foo".

Unicode character properties

Since PHP 4.4.0 and 5.1.0, three additional escape sequences to match generic character types are available when *UTF-8 mode* is selected. They are:

`\p{xx}`
a character with the *xx* property

`\P{xx}`
a character without the *xx* property

`\X`
an extended Unicode sequence

The property names represented by *xx* above are limited to the Unicode general category properties. Each character has exactly one such property, specified by a two-letter abbreviation. For compatibility with Perl, negation can be specified by including a circumflex between the opening brace and the property name. For example, `\p{^Lu}` is the same as `\P{Lu}`.

If only one letter is specified with `\p` or `\P`, it includes all the properties that start with that letter. In this case, in the absence of negation, the curly brackets in the escape sequence are optional; these two examples have the same effect:

`\p{L}`
`\pL`

Table 245. Supported property codes

| | |
|-----------|-------------------|
| <i>C</i> | Other |
| <i>Cc</i> | Control |
| <i>Cf</i> | Format |
| <i>Cn</i> | Unassigned |
| <i>Co</i> | Private use |
| <i>Cs</i> | Surrogate |
| <i>L</i> | Letter |
| <i>Ll</i> | Lower case letter |
| <i>Lm</i> | Modifier letter |
| <i>Lo</i> | Other letter |
| <i>Lt</i> | Title case letter |
| <i>Lu</i> | Upper case letter |
| <i>M</i> | Mark |
| <i>Mc</i> | Spacing mark |
| <i>Me</i> | Enclosing mark |
| <i>Mn</i> | Non-spacing mark |
| <i>N</i> | Number |

| | |
|-----------|-----------------------|
| <i>Nd</i> | Decimal number |
| <i>Nl</i> | Letter number |
| <i>No</i> | Other number |
| <i>P</i> | Punctuation |
| <i>Pc</i> | Connector punctuation |
| <i>Pd</i> | Dash punctuation |
| <i>Pe</i> | Close punctuation |
| <i>Pf</i> | Final punctuation |
| <i>Pi</i> | Initial punctuation |
| <i>Po</i> | Other punctuation |
| <i>Ps</i> | Open punctuation |
| <i>S</i> | Symbol |
| <i>Sc</i> | Currency symbol |
| <i>Sk</i> | Modifier symbol |
| <i>Sm</i> | Mathematical symbol |
| <i>So</i> | Other symbol |
| <i>Z</i> | Separator |
| <i>Zl</i> | Line separator |
| <i>Zp</i> | Paragraph separator |
| <i>Zs</i> | Space separator |

Extended properties such as "Greek" or "InMusicalSymbols" are not supported by PCRE.

Specifying caseless matching does not affect these escape sequences. For example, `\p{Lu}` always matches only upper case letters.

The `\X` escape matches any number of Unicode characters that form an extended Unicode sequence. `\X` is equivalent to `(?>\PM\pM*)`.

That is, it matches a character without the "mark" property, followed by zero or more characters with the "mark" property, and treats the sequence as an atomic group (see below). Characters with the "mark" property are typically accents that affect the preceding character.

Matching characters by Unicode property is not fast, because PCRE has to search a structure that contains data for over fifteen thousand characters. That is why the traditional escape sequences such as `\d` and `\w` do not use Unicode properties in PCRE.

Circumflex and dollar

Outside a character class, in the default matching mode, the circumflex character is an assertion which is true only if the current matching point is at the start of the subject string. Inside a character class, circumflex has an entirely different meaning (see below).

Circumflex need not be the first character of the pattern if a number of alternatives are involved, but it should be the first thing in each alternative in which it appears if the pattern is ever to match that branch. If all possible alternatives start with a circumflex, that is, if the pattern is constrained to match only at the start of the subject, it is said to be an "anchored" pattern. (There are also other constructs that can cause a pattern to be anchored.)

A dollar character is an assertion which is **TRUE** only if the current matching point is at the end of the subject string, or immediately before a newline character that is the last character in the string (by default). Dollar need not be the last character of the pattern if a number of alternatives are involved, but it should be the last item in any branch in which it appears. Dollar has no special meaning in a character class.

The meaning of dollar can be changed so that it matches only at the very end of the string, by setting the [PCRE DOLLAR ENDONLY](#) option at compile or matching time. This does not affect the `\Z` assertion.

The meanings of the circumflex and dollar characters are changed if the [PCRE MULTILINE](#) option is set. When this is the case, they match immediately after and immediately before an internal `"\n"` character, respectively, in addition to matching at the start and end of the subject string. For example, the pattern `/^abc$/` matches the subject string `"def\nabc"` in multiline mode, but not otherwise. Consequently, patterns that are anchored in single line mode because all branches start with `"^"` are not anchored in multiline mode. The [PCRE DOLLAR ENDONLY](#) option is ignored if [PCRE MULTILINE](#) is set.

Note that the sequences `\A`, `\Z`, and `\z` can be used to match the start and end of the subject in both modes, and if all branches of a pattern start with `\A` is it always anchored, whether [PCRE MULTILINE](#) is set or not.

Full stop

Outside a character class, a dot in the pattern matches any one character in the subject, including a non-printing character, but not (by default) newline. If the [PCRE DOTALL](#) option is set, then dots match newlines as well. The handling of dot is entirely independent of the handling of circumflex and dollar, the only relationship being that they both involve newline characters. Dot has no special meaning in a character class.

`\C` can be used to match single byte. It makes sense in *UTF-8 mode* where full stop matches the whole character which can consist of multiple bytes.

Square brackets

An opening square bracket introduces a character class, terminated by a closing square bracket. A closing square bracket on its own is not special. If a closing square bracket is required as a member of the class, it should be the first data character in the class (after an initial circumflex, if present) or escaped with a backslash.

A character class matches a single character in the subject; the character must be in the set of characters defined by the class, unless the first character in the class is a circumflex, in which case the subject character must not be in the set defined by the class. If a circumflex is actually required as a member of the class, ensure it is not the first character, or escape it with a backslash.

For example, the character class `[aeiou]` matches any lower case vowel, while `[^aeiou]` matches any character that is not a lower case vowel. Note that a circumflex is just a convenient notation for specifying the characters which are in the class by enumerating those that are not. It is not an assertion: it still consumes a character from the subject string, and fails if the current pointer is at the end of the string.

When caseless matching is set, any letters in a class represent both their upper case and lower case versions, so for example, a caseless `[aeiou]` matches "A" as well as "a", and a caseless `[^aeiou]` does not match "A", whereas a careful version would.

The newline character is never treated in any special way in character classes, whatever the setting of the [PCRE_DOTALL](#) or [PCRE_MULTILINE](#) options is. A class such as `[^a]` will always match a newline.

The minus (hyphen) character can be used to specify a range of characters in a character class. For example, `[d-m]` matches any letter between d and m, inclusive. If a minus character is required in a class, it must be escaped with a backslash or appear in a position where it cannot be interpreted as indicating a range, typically as the first or last character in the class.

It is not possible to have the literal character "]" as the end character of a range. A pattern such as `[W-]46]` is interpreted as a class of two characters ("W" and "-") followed by a literal string "46]", so it would match "W46]" or "-46]". However, if the "]" is escaped with a backslash it is interpreted as the end of range, so `[W-]46]` is interpreted as a single class containing a range followed by two separate characters. The octal or hexadecimal representation of "]" can also be used to end a range.

Ranges operate in ASCII collating sequence. They can also be used for characters specified numerically, for example `[000-037]`. If a range that includes letters is used when caseless matching is set, it matches the letters in either case. For example, `[W-c]` is equivalent to `[[]^_`wxyzabc]`, matched caselessly, and if character tables for the "fr" locale are in use, `[xc8-xcb]` matches accented E characters in both cases.

The character types `\d`, `\D`, `\s`, `\S`, `\w`, and `\W` may also appear in a character class, and add the characters that they match to the class. For example, `[\dABCDEF]` matches any hexadecimal digit. A circumflex can conveniently be used with the upper case character types to specify a more restricted set of characters than the matching lower case type. For example, the class `^[^W_]` matches any letter or digit, but not underscore.

All non-alphanumeric characters other than `\`, `-`, `^` (at the start) and the terminating `]` are non-special in character classes, but it does no harm if they are escaped.

Vertical bar

Vertical bar characters are used to separate alternative patterns. For example, the pattern *gilbert/sullivan* matches either "gilbert" or "sullivan". Any number of alternatives may appear, and an empty alternative is permitted (matching the empty string). The matching process tries each alternative in turn, from left to right, and the first one that succeeds is used. If the alternatives are within a subpattern (defined below), "succeeds" means matching the rest of the main pattern as well as the alternative in the subpattern.

Internal option setting

The settings of [PCRE_CASELESS](#), [PCRE_MULTILINE](#), [PCRE_DOTALL](#), [PCRE_UNGREEDY](#), [PCRE_EXTRA](#), [PCRE_EXTENDED](#) and `PCRE_DUPNAMES` can be changed from within the pattern

by a sequence of Perl option letters enclosed between "(" and ")". The option letters are:

Table 246. Internal option letters

| | |
|----------|--|
| <i>i</i> | for PCRE_CASELESS |
| <i>m</i> | for PCRE_MULTILINE |
| <i>s</i> | for PCRE_DOTALL |
| <i>x</i> | for PCRE_EXTENDED |
| <i>U</i> | for PCRE_UNGREEDY |
| <i>X</i> | for PCRE_EXTRA |
| <i>J</i> | for PCRE_INFO_JCHANGED |

For example, `(?im)` sets caseless, multiline matching. It is also possible to unset these options by preceding the letter with a hyphen, and a combined setting and unsetting such as `(?im-sx)`, which sets [PCRE_CASELESS](#) and [PCRE_MULTILINE](#) while unsetting [PCRE_DOTALL](#) and [PCRE_EXTENDED](#), is also permitted. If a letter appears both before and after the hyphen, the option is unset.

When an option change occurs at top level (that is, not inside subpattern parentheses), the change applies to the remainder of the pattern that follows. So `/ab(?i)c/` matches only "abc" and "abC". This behaviour has been changed in PCRE 4.0, which is bundled since PHP 4.3.3. Before those versions, `/ab(?i)c/` would perform as `/abc/i` (e.g. matching "ABC" and "aBc").

If an option change occurs inside a subpattern, the effect is different. This is a change of behaviour in Perl 5.005. An option change inside a subpattern affects only that part of the subpattern that follows it, so `(a(?i)b)c` matches abc and aBc and no other strings (assuming [PCRE_CASELESS](#) is not used). By this means, options can be made to have different settings in different parts of the pattern. Any changes made in one alternative do carry on into subsequent branches within the same subpattern. For example, `(a(?i)b/c)` matches "ab", "aB", "c", and "C", even though when matching "C" the first branch is abandoned before the option setting. This is because the effects of option settings happen at compile time. There would be some very weird behaviour otherwise.

The PCRE-specific options [PCRE_UNGREEDY](#) and [PCRE_EXTRA](#) can be changed in the same way as the Perl-compatible options by using the characters U and X respectively. The `(?X)` flag setting is special in that it must always occur earlier in the pattern than any of the additional features it turns on, even when it is at top level. It is best put at the start.

Subpatterns

Subpatterns are delimited by parentheses (round brackets), which can be nested. Marking part of a pattern as a subpattern does two things:

1. It localizes a set of alternatives. For example, the pattern `cat(aract|erpillar|)` matches one of the words "cat", "cataract", or "caterpillar". Without the parentheses, it would match "cataract", "erpillar" or the empty string.
2. It sets up the subpattern as a capturing subpattern (as defined above). When the whole pattern matches, that portion of the subject string that matched the subpattern is passed back to the caller via the

ovector argument of **pcre_exec()**. Opening parentheses are counted from left to right (starting from 1) to obtain the numbers of the capturing subpatterns.

For example, if the string "the red king" is matched against the pattern *the ((red/white) (king/queen))* the captured substrings are "red king", "red", and "king", and are numbered 1, 2, and 3.

The fact that plain parentheses fulfil two functions is not always helpful. There are often times when a grouping subpattern is required without a capturing requirement. If an opening parenthesis is followed by "?:", the subpattern does not do any capturing, and is not counted when computing the number of any subsequent capturing subpatterns. For example, if the string "the white queen" is matched against the pattern *the ((?:red/white) (king/queen))* the captured substrings are "white queen" and "queen", and are numbered 1 and 2. The maximum number of captured substrings is 99, and the maximum number of all subpatterns, both capturing and non-capturing, is 200.

As a convenient shorthand, if any option settings are required at the start of a non-capturing subpattern, the option letters may appear between the "?" and the ":". Thus the two patterns

```
(?i:saturday|sunday)
(?:(?i)saturday|sunday)
```

match exactly the same set of strings. Because alternative branches are tried from left to right, and options are not reset until the end of the subpattern is reached, an option setting in one branch does affect subsequent branches, so the above patterns match "SUNDAY" as well as "Saturday".

It is possible to name the subpattern with *(?P<name>pattern)* since PHP 4.3.3. Array with matches will contain the match indexed by the string alongside the match indexed by a number, then.

Repetition

Repetition is specified by quantifiers, which can follow any of the following items:

- a single character, possibly escaped
- the `.` metacharacter
- a character class
- a back reference (see next section)
- a parenthesized subpattern (unless it is an assertion - see below)

The general repetition quantifier specifies a minimum and maximum number of permitted matches, by giving the two numbers in curly brackets (braces), separated by a comma. The numbers must be less than 65536, and the first must be less than or equal to the second. For example: *z{2,4}* matches "zz", "zzz", or "zzzz". A closing brace on its own is not a special character. If the second number is omitted, but the comma is present, there is no upper limit; if the second number and the comma are both omitted, the quantifier specifies an exact number of required matches. Thus *[aeiou]{3,}* matches at least 3 successive vowels, but may match many more, while *\d{8}* matches exactly 8 digits. An opening curly bracket that appears in a position where a quantifier is not allowed, or one that does not match the syntax of a quantifier, is taken as a literal character. For example, *{,6}* is not a quantifier, but a literal string of four characters.

The quantifier `{0}` is permitted, causing the expression to behave as if the previous item and the quantifier were not present.

For convenience (and historical compatibility) the three most common quantifiers have single-character abbreviations:

Table 247. Single-character quantifiers

| | |
|---|----------------------------------|
| * | equivalent to <code>{0,}</code> |
| + | equivalent to <code>{1,}</code> |
| ? | equivalent to <code>{0,1}</code> |

It is possible to construct infinite loops by following a subpattern that can match no characters with a quantifier that has no upper limit, for example: `(a?)*`

Earlier versions of Perl and PCRE used to give an error at compile time for such patterns. However, because there are cases where this can be useful, such patterns are now accepted, but if any repetition of the subpattern does in fact match no characters, the loop is forcibly broken.

By default, the quantifiers are "greedy", that is, they match as much as possible (up to the maximum number of permitted times), without causing the rest of the pattern to fail. The classic example of where this gives problems is in trying to match comments in C programs. These appear between the sequences `/*` and `*/` and within the sequence, individual `*` and `/` characters may appear. An attempt to match C comments by applying the pattern `/*.**/` to the string `/* first comment */ not comment /* second comment */` fails, because it matches the entire string due to the greediness of the `.*` item.

However, if a quantifier is followed by a question mark, then it ceases to be greedy, and instead matches the minimum number of times possible, so the pattern `/*.*?*/` does the right thing with the C comments. The meaning of the various quantifiers is not otherwise changed, just the preferred number of matches. Do not confuse this use of question mark with its use as a quantifier in its own right. Because it has two uses, it can sometimes appear doubled, as in `\d??\d` which matches one digit by preference, but can match two if that is the only way the rest of the pattern matches.

If the [PCRE_UNGREEDY](#) option is set (an option which is not available in Perl) then the quantifiers are not greedy by default, but individual ones can be made greedy by following them with a question mark. In other words, it inverts the default behaviour.

Quantifiers followed by `+` are "possessive". They eat as many characters as possible and don't return to match the rest of the pattern. Thus `.*abc` matches "aabc" but `.*+abc` doesn't because `.*+` eats the whole string. Possessive quantifiers can be used to speed up processing since PHP 4.3.3.

When a parenthesized subpattern is quantified with a minimum repeat count that is greater than 1 or with a limited maximum, more store is required for the compiled pattern, in proportion to the size of the minimum or maximum.

If a pattern starts with `.*` or `{0,}` and the [PCRE_DOTALL](#) option (equivalent to Perl's `/s`) is set, thus allowing the `.` to match newlines, then the pattern is implicitly anchored, because whatever follows will be tried against every character position in the subject string, so there is no point in retrying the overall match at any position after the first. PCRE treats such a pattern as though it were preceded by `\A`. In

cases where it is known that the subject string contains no newlines, it is worth setting [PCRE_DOTALL](#) when the pattern begins with `.*` in order to obtain this optimization, or alternatively using `^` to indicate anchoring explicitly.

When a capturing subpattern is repeated, the value captured is the substring that matched the final iteration. For example, after `(tweedle[dume]{3}\s*)+` has matched "tweedledum tweedledee" the value of the captured substring is "tweedledee". However, if there are nested capturing subpatterns, the corresponding captured values may have been set in previous iterations. For example, after `/(a/(b))+/` matches "aba" the value of the second captured substring is "b".

Back references

Outside a character class, a backslash followed by a digit greater than 0 (and possibly further digits) is a back reference to a capturing subpattern earlier (i.e. to its left) in the pattern, provided there have been that many previous capturing left parentheses.

However, if the decimal number following the backslash is less than 10, it is always taken as a back reference, and causes an error only if there are not that many capturing left parentheses in the entire pattern. In other words, the parentheses that are referenced need not be to the left of the reference for numbers less than 10. See the section entitled "Backslash" above for further details of the handling of digits following a backslash.

A back reference matches whatever actually matched the capturing subpattern in the current subject string, rather than anything matching the subpattern itself. So the pattern `(sens/respons)e and \libility` matches "sense and sensibility" and "response and responsibility", but not "sense and responsibility". If careful matching is in force at the time of the back reference, then the case of letters is relevant. For example, `((?i)rah)\s+\I` matches "rah rah" and "RAH RAH", but not "RAH rah", even though the original capturing subpattern is matched caselessly.

There may be more than one back reference to the same subpattern. If a subpattern has not actually been used in a particular match, then any back references to it always fail. For example, the pattern `(a/(bc))\2` always fails if it starts to match "a" rather than "bc". Because there may be up to 99 back references, all digits following the backslash are taken as part of a potential back reference number. If the pattern continues with a digit character, then some delimiter must be used to terminate the back reference. If the [PCRE_EXTENDED](#) option is set, this can be whitespace. Otherwise an empty comment can be used.

A back reference that occurs inside the parentheses to which it refers fails when the subpattern is first used, so, for example, `(a\1)` never matches. However, such references can be useful inside repeated subpatterns. For example, the pattern `(a/b\I)+` matches any number of "a"s and also "aba", "ababaa" etc. At each iteration of the subpattern, the back reference matches the character string corresponding to the previous iteration. In order for this to work, the pattern must be such that the first iteration does not need to match the back reference. This can be done using alternation, as in the example above, or by a quantifier with a minimum of zero.

Back references to the named subpatterns can be achieved by `(?P=name)` or, since PHP 5.2.4, also by `\k<name>`, `\k'name'`, `\k{name}` or `\g{name}`.

Assertions

An assertion is a test on the characters following or preceding the current matching point that does not

actually consume any characters. The simple assertions coded as `\b`, `\B`, `\A`, `\Z`, `\z`, `^` and `$` are described above. More complicated assertions are coded as subpatterns. There are two kinds: those that look ahead of the current position in the subject string, and those that look behind it.

An assertion subpattern is matched in the normal way, except that it does not cause the current matching position to be changed. Lookahead assertions start with `(?=` for positive assertions and `(?!` for negative assertions. For example, `\w+(?=;)` matches a word followed by a semicolon, but does not include the semicolon in the match, and `foo(?!bar)` matches any occurrence of "foo" that is not followed by "bar". Note that the apparently similar pattern `(?!foo)bar` does not find an occurrence of "bar" that is preceded by something other than "foo"; it finds any occurrence of "bar" whatsoever, because the assertion `(?!foo)` is always **TRUE** when the next three characters are "bar". A lookbehind assertion is needed to achieve this effect.

Lookbehind assertions start with `(?<=` for positive assertions and `(?<!` for negative assertions. For example, `(?<!/foo)bar` does find an occurrence of "bar" that is not preceded by "foo". The contents of a lookbehind assertion are restricted such that all the strings it matches must have a fixed length. However, if there are several alternatives, they do not all have to have the same fixed length. Thus `(?<=bullock/donkey)` is permitted, but `(?<!/dogs?/cats?)` causes an error at compile time. Branches that match different length strings are permitted only at the top level of a lookbehind assertion. This is an extension compared with Perl 5.005, which requires all branches to match the same length of string. An assertion such as `(?<=ab(c/de))` is not permitted, because its single top-level branch can match two different lengths, but it is acceptable if rewritten to use two top-level branches: `(?<=abc/abde)`. The implementation of lookbehind assertions is, for each alternative, to temporarily move the current position back by the fixed width and then try to match. If there are insufficient characters before the current position, the match is deemed to fail. Lookbehinds in conjunction with once-only subpatterns can be particularly useful for matching at the ends of strings; an example is given at the end of the section on once-only subpatterns.

Several assertions (of any sort) may occur in succession. For example, `(?<=\d{3})(?<!999)foo` matches "foo" preceded by three digits that are not "999". Notice that each of the assertions is applied independently at the same point in the subject string. First there is a check that the previous three characters are all digits, then there is a check that the same three characters are not "999". This pattern does not match "foo" preceded by six characters, the first of which are digits and the last three of which are not "999". For example, it doesn't match "123abcfoo". A pattern to do that is `(?<=\d{3}...)(?<!999)foo`.

This time the first assertion looks at the preceding six characters, checking that the first three are digits, and then the second assertion checks that the preceding three characters are not "999".

Assertions can be nested in any combination. For example, `(?<=(?<!/foo)bar)baz` matches an occurrence of "baz" that is preceded by "bar" which in turn is not preceded by "foo", while `(?<=\d{3}...)(?<!999)foo` is another pattern which matches "foo" preceded by three digits and any three characters that are not "999".

Assertion subpatterns are not capturing subpatterns, and may not be repeated, because it makes no sense to assert the same thing several times. If any kind of assertion contains capturing subpatterns within it, these are counted for the purposes of numbering the capturing subpatterns in the whole pattern. However, substring capturing is carried out only for positive assertions, because it does not make sense for negative assertions.

Assertions count towards the maximum of 200 parenthesized subpatterns.

Once-only subpatterns

With both maximizing and minimizing repetition, failure of what follows normally causes the repeated item to be re-evaluated to see if a different number of repeats allows the rest of the pattern to match. Sometimes it is useful to prevent this, either to change the nature of the match, or to cause it fail earlier than it otherwise might, when the author of the pattern knows there is no point in carrying on.

Consider, for example, the pattern `\d+foo` when applied to the subject line *123456bar*

After matching all 6 digits and then failing to match "foo", the normal action of the matcher is to try again with only 5 digits matching the `\d+` item, and then with 4, and so on, before ultimately failing. Once-only subpatterns provide the means for specifying that once a portion of the pattern has matched, it is not to be re-evaluated in this way, so the matcher would give up immediately on failing to match "foo" the first time. The notation is another kind of special parenthesis, starting with `(?>` as in this example: `(?>\d+)bar`

This kind of parenthesis "locks up" the part of the pattern it contains once it has matched, and a failure further into the pattern is prevented from backtracking into it. Backtracking past it to previous items, however, works as normal.

An alternative description is that a subpattern of this type matches the string of characters that an identical standalone pattern would match, if anchored at the current point in the subject string.

Once-only subpatterns are not capturing subpatterns. Simple cases such as the above example can be thought of as a maximizing repeat that must swallow everything it can. So, while both `\d+` and `\d+?` are prepared to adjust the number of digits they match in order to make the rest of the pattern match, `(?>\d+)` can only match an entire sequence of digits.

This construction can of course contain arbitrarily complicated subpatterns, and it can be nested.

Once-only subpatterns can be used in conjunction with look-behind assertions to specify efficient matching at the end of the subject string. Consider a simple pattern such as *abcd\$* when applied to a long string which does not match. Because matching proceeds from left to right, PCRE will look for each "a" in the subject and then see if what follows matches the rest of the pattern. If the pattern is specified as `^.*abcd$` then the initial `.*` matches the entire string at first, but when this fails (because there is no following "a"), it backtracks to match all but the last character, then all but the last two characters, and so on. Once again the search for "a" covers the entire string, from right to left, so we are no better off. However, if the pattern is written as `^(?>.*)(?<=abcd)` then there can be no backtracking for the `.*` item; it can match only the entire string. The subsequent lookbehind assertion does a single test on the last four characters. If it fails, the match fails immediately. For long strings, this approach makes a significant difference to the processing time.

When a pattern contains an unlimited repeat inside a subpattern that can itself be repeated an unlimited number of times, the use of a once-only subpattern is the only way to avoid some failing matches taking a very long time indeed. The pattern `(\D+<\d+>)*[!?]` matches an unlimited number of substrings that either consist of non-digits, or digits enclosed in `<>`, followed by either `!` or `?`. When it matches, it runs quickly. However, if it is applied to *aaa* it takes a long time before reporting failure. This is because the string can be divided between the two repeats in a large number of ways, and all have to be tried. (The example used `[!?]` rather than a single character at the end, because both PCRE and Perl have an optimization that allows for fast failure when

a single character is used. They remember the last single character that is required for a match, and fail early if it is not present in the string.) If the pattern is changed to `((?>\D+)/<\d+>)*[!/?]` sequences of non-digits cannot be broken, and failure happens quickly.

Conditional subpatterns

It is possible to cause the matching process to obey a subpattern conditionally or to choose between two alternative subpatterns, depending on the result of an assertion, or whether a previous capturing subpattern matched or not. The two possible forms of conditional subpattern are

```
(?(condition)yes-pattern)
(?(condition)yes-pattern|no-pattern)
```

If the condition is satisfied, the yes-pattern is used; otherwise the no-pattern (if present) is used. If there are more than two alternatives in the subpattern, a compile-time error occurs.

There are two kinds of condition. If the text between the parentheses consists of a sequence of digits, then the condition is satisfied if the capturing subpattern of that number has previously matched. Consider the following pattern, which contains non-significant white space to make it more readable (assume the [PCRE_EXTENDED](#) option) and to divide it into three parts for ease of discussion: `(\ () ? [^ ()] + (? (1) \))`

The first part matches an optional opening parenthesis, and if that character is present, sets it as the first captured substring. The second part matches one or more characters that are not parentheses. The third part is a conditional subpattern that tests whether the first set of parentheses matched or not. If they did, that is, if subject started with an opening parenthesis, the condition is **TRUE**, and so the yes-pattern is executed and a closing parenthesis is required. Otherwise, since no-pattern is not present, the subpattern matches nothing. In other words, this pattern matches a sequence of non-parentheses, optionally enclosed in parentheses.

If the condition is the string `(R)`, it is satisfied if a recursive call to the pattern or subpattern has been made. At "top level", the condition is false.

If the condition is not a sequence of digits or `(R)`, it must be an assertion. This may be a positive or negative lookahead or lookbehind assertion. Consider this pattern, again containing non-significant white space, and with the two alternatives on the second line:

```
(?(?=[^a-z]*[a-z])
\d{2}-[a-z]{3}-\d{2} | \d{2}-\d{2}-\d{2} )
```

The condition is a positive lookahead assertion that matches an optional sequence of non-letters followed by a letter. In other words, it tests for the presence of at least one letter in the subject. If a letter is found, the subject is matched against the first alternative; otherwise it is matched against the second. This pattern matches strings in one of the two forms `dd-aaa-dd` or `dd-dd-dd`, where `aaa` are letters and `dd` are digits.

Comments

The sequence `(?#` marks the start of a comment which continues up to the next closing parenthesis. Nested parentheses are not permitted. The characters that make up a comment play no part in the pattern matching at all.

If the [PCRE_EXTENDED](#) option is set, an unescaped `#` character outside a character class introduces a comment that continues up to the next newline character in the pattern.

Recursive patterns

Consider the problem of matching a string in parentheses, allowing for unlimited nested parentheses. Without the use of recursion, the best that can be done is to use a pattern that matches up to some fixed depth of nesting. It is not possible to handle an arbitrary nesting depth. Perl 5.6 has provided an experimental facility that allows regular expressions to recurse (among other things). The special item `(?R)` is provided for the specific case of recursion. This PCRE pattern solves the parentheses problem (assume the [PCRE_EXTENDED](#) option is set so that white space is ignored): `\(((?>[^()]+) / (?R))* \)`

First it matches an opening parenthesis. Then it matches any number of substrings which can either be a sequence of non-parentheses, or a recursive match of the pattern itself (i.e. a correctly parenthesized substring). Finally there is a closing parenthesis.

This particular example pattern contains nested unlimited repeats, and so the use of a once-only subpattern for matching strings of non-parentheses is important when applying the pattern to strings that do not match. For example, when it is applied to `(aaa)` it yields "no match" quickly. However, if a once-only subpattern is not used, the match runs for a very long time indeed because there are so many different ways the `+` and `*` repeats can carve up the subject, and all have to be tested before failure can be reported.

The values set for any capturing subpatterns are those from the outermost level of the recursion at which the subpattern value is set. If the pattern above is matched against `ab(cd)ef` the value for the capturing parentheses is "ef", which is the last value taken on at the top level. If additional parentheses are added, giving `\(((?>[^()]+) / (?R))* \)` then the string they capture is "ab(cd)ef", the contents of the top level parentheses. If there are more than 15 capturing parentheses in a pattern, PCRE has to obtain extra memory to store data during a recursion, which it does by using `pcre_malloc`, freeing it via `pcre_free` afterwards. If no memory can be obtained, it saves data for the first 15 capturing parentheses only, as there is no way to give an out-of-memory error from within a recursion.

Since PHP 4.3.3, `(?1)`, `(?2)` and so on can be used for recursive subpatterns too. It is also possible to use named subpatterns: `(?P>name)` or `(?P&name)`.

If the syntax for a recursive subpattern reference (either by number or by name) is used outside the parentheses to which it refers, it operates like a subroutine in a programming language. An earlier example pointed out that the pattern `(sens/respons)e and \libility` matches "sense and sensibility" and "response and responsibility", but not "sense and responsibility". If instead the pattern `(sens/respons)e and (?1)ibility` is used, it does match "sense and responsibility" as well as the other two strings. Such references must, however, follow the subpattern to which they refer.

The maximum length of a subject string is the largest positive number that an integer variable can hold.

However, PCRE uses recursion to handle subpatterns and indefinite repetition. This means that the available stack space may limit the size of a subject string that can be processed by certain patterns.

Performances

Certain items that may appear in patterns are more efficient than others. It is more efficient to use a character class like `[aeiou]` than a set of alternatives such as `(a|e|i|o|u)`. In general, the simplest construction that provides the required behaviour is usually the most efficient. Jeffrey Friedl's book contains a lot of discussion about optimizing regular expressions for efficient performance.

When a pattern begins with `.*` and the [PCRE_DOTALL](#) option is set, the pattern is implicitly anchored by PCRE, since it can match only at the start of a subject string. However, if [PCRE_DOTALL](#) is not set, PCRE cannot make this optimization, because the `.` metacharacter does not then match a newline, and if the subject string contains newlines, the pattern may match from the character immediately following one of them instead of from the very start. For example, the pattern `(.*) second` matches the subject `"first\nand second"` (where `\n` stands for a newline character) with the first captured substring being `"and"`. In order to do this, PCRE has to retry the match starting after every newline in the subject.

If you are using such a pattern with subject strings that do not contain newlines, the best performance is obtained by setting [PCRE_DOTALL](#), or starting the pattern with `^.*` to indicate explicit anchoring. That saves PCRE from having to scan along the subject looking for a newline to restart at.

Beware of patterns that contain nested indefinite repeats. These can take a long time to run when applied to a string that does not match. Consider the pattern fragment `(a+)*`

This can match `"aaaa"` in 33 different ways, and this number increases very rapidly as the string gets longer. (The `*` repeat can match 0, 1, 2, 3, or 4 times, and for each of those cases other than 0, the `+` repeats can match different numbers of times.) When the remainder of the pattern is such that the entire match is going to fail, PCRE has in principle to try every possible variation, and this can take an extremely long time.

An optimization catches some of the more simple cases such as `(a+)*b` where a literal character follows. Before embarking on the standard matching procedure, PCRE checks that there is a `"b"` later in the subject string, and if there is not, it fails the match immediately. However, when there is no following literal this optimization cannot be used. You can see the difference by comparing the behaviour of `(a+)*\d` with the pattern above. The former gives a failure almost instantly when applied to a whole line of `"a"` characters, whereas the latter takes an appreciable time with strings longer than about 20 characters.

[Prev](#)

Pattern Modifiers

[Up](#)[Home](#)[Next](#)

preg_grep