

Overview		Package	Class	Use	Tree	Deprecated	Index	Help
PREV CLASS	NEXT CLASS	FRAMES	NO FRAMES	All Classes	DETAIL: FIELD CONSTR METHOD			
CLASS Pattern								
java.util.regex	java.lang.Object			java.util.regex.Pattern				
All Implemented Interfaces:								
				Serializable				
<hr/>								
Java™ 2 Platform Std. Ed. v1.4.2								
javauil.regex								
<hr/>								
Character classes								
POSIX character classes (US-ASCII only)								
<hr/>								
Construct								
Matches								
Characters								
Classes for Unicode blocks and categories								
<hr/>								
http://java.sun.com/j2se/1.4.2/doc...								

Line terminators

A *line terminator* is a one- or two-character sequence that marks the end of a line of the input character sequence. The following are recognized as line terminators:

- A newline (line feed) character ('\\n'),
- A carriage-return character followed immediately by a newline character ('\\r\\n').
- A standalone carriage-return character ('\\r'),
- A next-line character ('\u0085'),
- A line-separator character ('\u2028'), or
- A paragraph-separator character ('\u2029').

If

`UNIX_LINES` mode is activated, then the only line terminators recognized are newline characters.

The regular expression `.` matches any character except a line terminator unless the `DOTALL` flag is specified.

By default, the regular expressions `^` and `$` ignore line terminators and only match at the beginning and the end, respectively, of the entire input sequence. If `MULTILINE` mode is activated then `^` matches at the beginning of input and after any line terminator except at the end of input. When in `MULTILINE` mode `$` matches just before a line terminator or the end of the input sequence.

Groups and capturing

Capturing groups are numbered by counting their opening parentheses from left to right. In the expression `((A)(B(C)))`, for example, there are four such groups:

```
1 ((A)(B(C)))
2 (A)
3 (B(C))
4 (C)
```

Group zero always stands for the entire expression.

Capturing groups are so named because, during a match, each subsequence of the input sequence that matches such a group is saved. The captured subsequence may be used later in the expression, via a back reference, and may also be retrieved from the matcher once the match operation is complete.

The captured input associated with a group is always the subsequence that the group most recently matched. If a group is evaluated a second time because of quantification then its previously-captured value, if any, will be retained if the second evaluation fails. Matching the string "aba" against the expression "(a(b)?)*", for example, leaves group two set to "ba". All captured input is discarded at the beginning of each match.

Groups beginning with `(?` are pure, *non-capturing* groups that do not capture text and do not count towards the group total.

Unicode support

This class follows

[Unicode Technical Report #18: Unicode Regular Expression Guidelines](#), implementing its second level of support though with a slightly different concrete syntax.

Unicode escape sequences such as '\u2014' in Java source code are processed as described in [23.3](#) of the Java Language Specification. Such escape sequences are also implemented directly by the regular-expression parser so that Unicode escapes can be used in expressions that are read from files or from the keyboard. Thus the strings "\u2014" and "\u2014" both match the same pattern, which matches the character with hexadecimal value 0x2014.

Unicode blocks and categories are written with the `\p` and `\P` constructs as in Perl. `\p{prop}` matches if the input has the property `prop`, while `\P{prop}` does not match if the input has that property. Blocks are specified with the prefix `\n`, as in `\nBasic_Latin`. Categories may be specified with the optional prefix `\s`: Both `\p{\L}` and `\p{\ISL}` denote the category of Unicode letters. Blocks and categories can be used both inside and outside of a character class.

The supported blocks and categories are those of

[The Unicode Standard, Version 3.0](#). The block names are those defined in Chapter 14 and in the file `Blocks-3.txt` of the [Unicode Character Database](#) except that the spaces are removed; "Basic Latin", for example, becomes "BasicLatin". The category names are those defined in table 4-5 of the Standard (p. 88), both normative and informative.

Comparison to Perl 5

Perl constructs not supported by this class:

- The conditional constructs `(?{X})` and `(?(condition) X|Y)`,
- The embedded code constructs `?{(code)}` and `?{?(code)}`,
- The embedded comment syntax `(?#comment)`, and
- The preprocessing operations `\U`, `\L`, and `\u`.

Constructs supported by this class but not by Perl:

- Possessive quantifiers, which greedily match as much as they can and do not back off, even when doing so would allow the overall match to succeed.
 - Character-class union and intersection as described [above](#).
- Notable differences from Perl:
- In Perl, `\1` through `\9` are always interpreted as back references; a backslash-escaped number greater than 9 is treated as a back reference if at least that many subexpressions exist, otherwise it is interpreted, if possible, as an octal escape. In this class octal escapes must always begin with a zero. In this class, `\1` through `\9` are always interpreted as back references, and a larger number is accepted as a back reference if at least that many subexpressions exist at that point in the regular expression, otherwise the parser will drop digits until the number is smaller or equal to the existing number of groups or it is one digit.
 - Perl uses the `g` flag to request a match that resumes where the last match left off. This functionality is provided implicitly by the `Matcher` class: Repeated invocations of the `Find` method will resume where the last match left off, unless the matcher is reset.

- In Perl, embedded flags at the top level of an expression affect the whole expression. In this class, embedded flags always take effect at the point at which they appear, whether they are at the top level or within a group; in the latter case, flags are restored at the end of the group just as in Perl.
- Perl is forgiving about malformed matching constructs, as in the expression `*a`, as well as dangling brackets, as in the expression `abc]`, and treats them as literals. This class also accepts dangling brackets but is strict about dangling metacharacters like `+`, `?` and `*`, and will throw a `PatternSyntaxException` if it encounters them.

For a more precise description of the behavior of regular expression constructs, please see [Mastering Regular Expressions, 2nd Edition](#). Jeffrey E. F. Friedl, O'Reilly and Associates, 2002.

Since:

1.4

See Also:

`String.split(String, int)`, `String.split(String, SerializedForm)`

Field Summary

static int CANON_EQ	Enables canonical equivalence.
static int CASE_INSENSITIVE	Enables case-insensitive matching.
static int COMMENTS	Permits whitespace and comments in pattern.
static int DOTALL	Enables dotall mode.
static int MULTILINE	Enables multiline mode.
static int UNICODE_CASE	Enables Unicode-aware case folding.
static int UNIX_LINES	Enables Unix lines mode.

Method Summary

static Pattern compile(String, regex)	Compiles the given regular expression into a pattern.
static Pattern compile(String, regex, int flags)	Compiles the given regular expression into a pattern with the given flags.
int flags()	Returns this pattern's match flags.
Matcher matcher(CharSequence, input)	Creates a matcher that will match the given input against this pattern.
static boolean matches(String, CharSequence, input)	Compiles the given regular expression and attempts to match the given input against it.
String Pattern()	Returns the regular expression from which this pattern was compiled.

- In Perl, embedded flags always take effect at the point at which they appear, whether they are at the top level or within a group; in the latter case, flags are restored at the end of the group just as in Perl.
- Perl is forgiving about malformed matching constructs, as in the expression `*a`, as well as dangling brackets, as in the expression `abc]`, and treats them as literals. This class also accepts dangling brackets but is strict about dangling metacharacters like `+`, `?` and `*`, and will throw a `PatternSyntaxException` if it encounters them.

- For a more precise description of the behavior of regular expression constructs, please see [Mastering Regular Expressions, 2nd Edition](#). Jeffrey E. F. Friedl, O'Reilly and Associates, 2002.

Since:

1.4

See Also:

`String.split(String, int)`, `String.split(String, SerializedForm)`

Field Detail

UNIX_LINES

public static final int [UNIX_LINES](#)

Enables Unix lines mode.

In this mode, only the `\n` line terminator is recognized in the behavior of `,`, `*`, and `$`.

Unix lines mode can also be enabled via the embedded flag expression `(?d)`.

See Also:

[Constant Field Values](#)

CASE_INSENSITIVE

public static final int [CASE_INSENSITIVE](#)

Enables case-insensitive matching.

By default, case-insensitive matching assumes that only characters in the US-ASCII charset are being matched. Unicode-aware case-insensitive matching can be enabled by specifying the `UNICODE_CASE` flag in conjunction with this flag.

Case-insensitive matching can also be enabled via the embedded flag expression `(?i)`.

Specifying this flag may impose a slight performance penalty.

See Also:

[Constant Field Values](#)

COMMENTS

public static final int [COMMENTS](#)

Permits whitespace and comments in pattern.

In this mode, whitespace is ignored, and embedded comments starting with `#` are ignored until the end of a line.

Comments mode can also be enabled via the embedded flag expression (?xi).

See Also:

[Constant Field Values](#)

MULTILINE

```
public static final int MULTILINE
```

Enables multiline mode.

In multiline mode the expressions ^ and \$ match just after or just before, respectively, a line terminator or the end of the input sequence. By default these expressions only match at the beginning and the end of the entire input sequence.

Multiline mode can also be enabled via the embedded flag expression (?m).

See Also:

[Constant Field Values](#)

DOTALL

```
public static final int DOTALL
```

Enables dotall mode.

In dotall mode, the expression . matches any character, including a line terminator. By default this expression does not match line terminators.

Dotall mode can also be enabled via the embedded flag expression (?s). (The s is a mnemonic for "single-line" mode, which is what this is called in Perl.)

See Also:

[Constant Field Values](#)

UNICODE_CASE

```
public static final int UNICODE_CASE
```

Enables Unicode-aware case folding.

When this flag is specified then case-insensitive matching, when enabled by the [CASE_INSENSITIVE](#) flag, is done in a manner consistent with the Unicode Standard. By default, case-insensitive matching assumes that only characters in the US-ASCII charset are being matched.

Unicode-aware case folding can also be enabled via the embedded flag expression (?u).

Specifying this flag may impose a performance penalty.

See Also:

[Constant Field Values](#)

Comments mode can also be enabled via the embedded flag expression (?xi).

See Also:

[Constant Field Values](#)

CANON_EQ

```
public static final int CANON_EQ
```

Enables canonical equivalence.

When this flag is specified then two characters will be considered to match if, and only if, their full canonical decompositions match. The expression `\a\ud830\udc0A`, for example, will match the string " ? " when this flag is specified. By default, matching does not take canonical equivalence into account.

There is no embedded flag character for enabling canonical equivalence.

Specifying this flag may impose a performance penalty.

See Also:

[Constant Field Values](#)

Method Detail

compile

```
public static Pattern compile(String regex)
```

Compiles the given regular expression into a pattern.

Parameters:

regex - The expression to be compiled

Throws:

[PatternSyntaxException](#) - If the expression's syntax is invalid

compile

```
public static Pattern compile(String regex, int flags)
```

Compiles the given regular expression into a pattern with the given flags.

Parameters:

regex - The expression to be compiled

flags - Match flags, a bit mask that may include [CASE_INSENSITIVE](#), [MULTILINE](#), [DOTALL](#), [UNICODE_CASE](#), and [CANON_EQ](#)

Throws:

[IllegalArgumentException](#) - If bit values other than those corresponding to the defined match flags are set in flags
[PatternSyntaxException](#) - If the expression's syntax is invalid

pattern

```
public String pattern()
```

Returns the regular expression from which this pattern was compiled.

Returns:

The source of this pattern

split

```
public String[] split(CharSequence input,
                     int limit)
```

Splits the given input sequence around matches of this pattern.

The array returned by this method contains each substring of the input sequence that is terminated by another subsequence that matches this pattern or is terminated by the end of the input sequence. The substrings in the array are in the order in which they occur in the input. If this pattern does not match any subsequence of the input then the resulting array has just one element, namely the input sequence in string form.

The `limit` parameter controls the number of times the pattern is applied and therefore affects the length of the resulting array. If the `limit` *n* is greater than zero then the pattern will be applied at most *n* - 1 times; the array's length will be no greater than *n*, and the array's last entry will contain all input beyond the last matched delimiter. If *n* is non-positive then the pattern will be applied as many times as possible and the array can have any length. If *n* is zero then the pattern will be applied as many times as possible, the array can have any length, and trailing empty strings will be discarded.

The input "boo:and:foo", for example, yields the following results with these parameters:

<i>Regex</i>	<i>Limit</i>	<i>Result</i>
:	2	{ "boo", "and:foo" }
:	5	{ "boo", "and:", "foo" }
:	-2	{ "boo", "and:", "foo" }
o	5	{ "b", "", "and:f", "" }
o	-2	{ "b", "", "and:f", "", "" }
o	0	{ "b", "", "and:f" }

Parameters:

`input` - The character sequence to be split

`limit` - The result threshold, as described above

Returns:

The array of strings computed by splitting the input around matches of this pattern

split

```
public String[] split(CharSequence input)
```

Splits the given input sequence around matches of this pattern.

This method works as if by invoking the two-argument `split` method with the given input sequence and a limit argument of zero. Trailing empty strings are therefore not included in the resulting array.

The input "boo:and:foo", for example, yields the following results with these expressions:

<i>Regex</i>	<i>Result</i>
:	{ "boo", "and", "foo" }

```

o   { "b", "", "and:f" }

```

Parameters:

`input` - The character sequence to be split

Returns:

The array of strings computed by splitting the input around matches of this pattern

Overview	Package	Class	Use	Tree	Deprecated	Index	Help	Frames	No Frames	All Classes	Java™ 2 Platform Std. Ed. v1.4.2
PREV CLASS	NEXT CLASS							DETAIL	FIELD	CONSTR	METHOD
SUMMARY: NESTED FIELD CONSTR METHOD											

[Submit a bug or feature](#)
For further API reference and developer documentation, see [java 2 SDK SE Developer Documentation](#). That documentation contains more detailed, developer-targeted descriptions, with conceptual overviews, definitions of terms, workarounds, and working code examples.

Copyright 2003 Sun Microsystems, Inc. All rights reserved. Use is subject to [license terms](#). Also see the [documentation redistribution policy](#).

<http://java.sun.com/j2se/1.4.2/doc/>