

Construct	Matches	Characters
<i>x</i>	The character <i>x</i>	
<code>\\</code>	The backslash character	
<code>\0n</code>	The character with octal value <i>0n</i> ( $0 \leq n \leq 7$ )	
<code>\0nn</code>	The character with octal value <i>0nn</i> ( $0 \leq n \leq 7$ )	
<code>\0mnn</code>	The character with octal value <i>0mnn</i> ( $0 \leq m \leq 3, 0 \leq n \leq 7$ )	
<code>\xhh</code>	The character with hexadecimal value <i>0xhh</i>	
<code>\uhhhh</code>	The character with hexadecimal value <i>0xhhhh</i>	
<code>\x{h...h}</code>	The character with hexadecimal value <i>0xh...h</i> ( <code>Character.MIN_CODE_POINT</code>	
<code>\t</code>	The tab character ( <code>'\u0009'</code> )	
<code>\n</code>	The newline (line feed) character ( <code>'\u000A'</code> )	
<code>\r</code>	The carriage-return character ( <code>'\u000D'</code> )	
<code>\f</code>	The form-feed character ( <code>'\u000C'</code> )	
<code>\a</code>	The alert (bell) character ( <code>'\u0007'</code> )	
<code>\e</code>	The escape character ( <code>'\u001B'</code> )	
<code>\cx</code>	The control character corresponding to <i>x</i>	
<b>Character classes</b>		
<code>[abc]</code>	a, b, or c (simple class)	
<code>[^abc]</code>	Any character except a, b, or c (negation)	
<code>[a-zA-Z]</code>	a through z or A through Z, inclusive (range)	
<code>[a-d[m-p]]</code>	a through d, or m through p: <code>[a-dm-p]</code> (union)	
<code>[a-z&amp;&amp;[def]]</code>	d, e, or f (intersection)	
<code>[a-z&amp;&amp;[^bc]]</code>	a through z, except for b and c: <code>[ad-z]</code> (subtraction)	
<code>[a-z&amp;&amp;[^m-p]]</code>	a through z, and not m through p: <code>[a-lq-z]</code> (subtraction)	
<b>Predefined character classes</b>		
<code>.</code>	Any character (may or may not match <a href="#">line terminators</a> )	
<code>\d</code>	A digit: <code>[0-9]</code>	
<code>\D</code>	A non-digit: <code>[^0-9]</code>	
<code>\s</code>	A whitespace character: <code>[ \t\n\x0B\f\r]</code>	
<code>\S</code>	A non-whitespace character: <code>[^\s]</code>	

<code>\w</code>	A word character: <code>[a-zA-Z_0-9]</code>
<code>\W</code>	A non-word character: <code>[^\w]</code>

### POSIX character classes (US-ASCII)

<code>\p{Lower}</code>	A lower-case alphabetic character: <code>[a-z]</code>
<code>\p{Upper}</code>	An upper-case alphabetic character: <code>[A-Z]</code>
<code>\p{ASCII}</code>	All ASCII: <code>[\x00-\x7F]</code>
<code>\p{Alpha}</code>	An alphabetic character: <code>[\p{Lower}\p{Upper}]</code>
<code>\p{Digit}</code>	A decimal digit: <code>[0-9]</code>
<code>\p{Alnum}</code>	An alphanumeric character: <code>[\p{Alpha}\p{Digit}]</code>
<code>\p{Punct}</code>	Punctuation: One of <code>!"#\$%&amp;'()*+,-./:;&lt;=&gt;?@[ \]^_`</code>
<code>\p{Graph}</code>	A visible character: <code>[\p{Alnum}\p{Punct}]</code>
<code>\p{Print}</code>	A printable character: <code>[\p{Graph}\x20]</code>
<code>\p{Blank}</code>	A space or a tab: <code>[\t]</code>
<code>\p{Cntrl}</code>	A control character: <code>[\x00-\x1F\x7F]</code>
<code>\p{XDigit}</code>	A hexadecimal digit: <code>[0-9a-fA-F]</code>
<code>\p{Space}</code>	A whitespace character: <code>[\t\n\x0B\f\r]</code>

### java.lang.Character classes (simple [java char](#))

<code>\p{javaLowerCase}</code>	Equivalent to <code>java.lang.Character.isLowerCase()</code>
<code>\p{javaUpperCase}</code>	Equivalent to <code>java.lang.Character.isUpperCase()</code>
<code>\p{javaWhitespace}</code>	Equivalent to <code>java.lang.Character.isWhitespace()</code>
<code>\p{javaMirrored}</code>	Equivalent to <code>java.lang.Character.isMirrored()</code>

### Classes for Unicode scripts, blocks, categories and properties

<code>\p{IsLatin}</code>	A Latin script character ( <a href="#">script</a> )
<code>\p{InGreek}</code>	A character in the Greek block ( <a href="#">block</a> )
<code>\p{Lu}</code>	An uppercase letter ( <a href="#">category</a> )
<code>\p{IsAlphabetic}</code>	An alphabetic character ( <a href="#">binary property</a> )
<code>\p{Sc}</code>	A currency symbol
<code>\P{InGreek}</code>	Any character except one in the Greek block (negation)
<code>[ \p{L} &amp;&amp; [ ^ \p{Lu} ] ]</code>	Any letter except an uppercase letter (subtraction)

### Boundary matchers

<code>^</code>	The beginning of a line
----------------	-------------------------

\$	The end of a line
\b	A word boundary
\B	A non-word boundary
\A	The beginning of the input
\G	The end of the previous match
\Z	The end of the input but for the final <u>terminator</u> , if any
\z	The end of the input

### Greedy quantifiers

$X?$	$X$ , once or not at all
$X^*$	$X$ , zero or more times
$X^+$	$X$ , one or more times
$X\{n\}$	$X$ , exactly $n$ times
$X\{n, \}$	$X$ , at least $n$ times
$X\{n, m\}$	$X$ , at least $n$ but not more than $m$ times

### Reluctant quantifiers

$X??$	$X$ , once or not at all
$X^*?$	$X$ , zero or more times
$X^+?$	$X$ , one or more times
$X\{n\}?$	$X$ , exactly $n$ times
$X\{n, \}?$	$X$ , at least $n$ times
$X\{n, m\}?$	$X$ , at least $n$ but not more than $m$ times

### Possessive quantifiers

$X^?+$	$X$ , once or not at all
$X^*+$	$X$ , zero or more times
$X^{++}$	$X$ , one or more times
$X\{n\}^+$	$X$ , exactly $n$ times
$X\{n, \}^+$	$X$ , at least $n$ times
$X\{n, m\}^+$	$X$ , at least $n$ but not more than $m$ times

### Logical operators

$XY$	$X$ followed by $Y$
$X Y$	Either $X$ or $Y$

(X) X, as a [capturing group](#)

## Back references

\n Whatever the  $n^{\text{th}}$  [capturing group](#) matched

\k<name> Whatever the [named-capturing group](#) "name" matched

## Quotation

\ Nothing, but quotes the following character

\Q Nothing, but quotes all characters until \E

\E Nothing, but ends quoting started by \Q

## Special constructs (named-capturing and non-capturing)

(?<[name](#)>X) X, as a named-capturing group

(?:X) X, as a non-capturing group

(?idmsuxU-idmsuxU) Nothing, but turns match flags [i d m s u x U](#) on - off

(?idmsux-idmsux:X) X, as a [non-capturing group](#) with the given flags [i d m s u x](#) on - off

(?=X) X, via zero-width positive lookahead

(?!X) X, via zero-width negative lookahead

(?<=X) X, via zero-width positive lookbehind

(?<!X) X, via zero-width negative lookbehind

(?>X) X, as an independent, non-capturing group

---

## Backslashes, escapes, and quoting

The backslash character ( '\ ' ) serves to introduce escaped constructs, as defined in the table above, as well as to quote characters that otherwise would be interpreted as unescaped constructs. Thus the expression \\ matches a single backslash and \{ matches a left brace.

It is an error to use a backslash prior to any alphabetic character that does not denote an escaped construct; these are reserved for future extensions to the regular-expression language. A backslash may be used prior to a non-alphabetic character regardless of whether that character is part of an unescaped construct.

Backslashes within string literals in Java source code are interpreted as required by *The Java™ Language Specification* as either Unicode escapes (section 3.3) or other character escapes (section 3.10.6) It is therefore necessary to double backslashes in string literals that represent regular expressions to protect them from interpretation by the Java bytecode compiler. The string literal "\\b", for example, matches a single backspace character when interpreted as a regular expression, while "\\b" matches a word boundary. The string literal "\\ (hello\\) " is illegal and leads to a compile-time error; in order to match the string (hello) the string literal "\\ (hello\\) " must be used.

## Character Classes

Character classes may appear within other character classes, and may be composed by the union operator (implicit) and the intersection operator (& &). The union operator denotes a class that contains every character

that is in at least one of its operand classes. The intersection operator denotes a class that contains every character that is in both of its operand classes.

The precedence of character-class operators is as follows, from highest to lowest:

<b>1</b>	Literal escape	<code>\x</code>
<b>2</b>	Grouping	<code>[...]</code>
<b>3</b>	Range	<code>a-z</code>
<b>4</b>	Union	<code>[a-e][i-u]</code>
<b>5</b>	Intersection	<code>[a-z&amp;&amp;[aeiou]]</code>

Note that a different set of metacharacters are in effect inside a character class than outside a character class. For instance, the regular expression `.` loses its special meaning inside a character class, while the expression `-` becomes a range forming metacharacter.

## 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

### Group number

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:

<b>1</b>	<code>((A)(B(C)))</code>
<b>2</b>	<code>(A)</code>
<b>3</b>	<code>(B(C))</code>
<b>4</b>	<code>(C)</code>

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.

## Group name

A capturing group can also be assigned a "name", a `named-capturing group`, and then be back-referenced later by the "name". Group names are composed of the following characters. The first character must be a letter.

- The uppercase letters 'A' through 'Z' ('`\u0041`' through '`\u005a`'),
- The lowercase letters 'a' through 'z' ('`\u0061`' through '`\u007a`'),
- The digits '0' through '9' ('`\u0030`' through '`\u0039`'),

A `named-capturing group` is still numbered as described in [Group number](#).

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 "b". All captured input is discarded at the beginning of each match.

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

## Unicode support

This class is in conformance with Level 1 of [Unicode Technical Standard #18: Unicode Regular Expression](#), plus RL2.1 Canonical Equivalents.

**Unicode escape sequences** such as `\u2014` in Java source code are processed as described in section 3.3 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"`, while not equal, compile into the same pattern, which matches the character with hexadecimal value `0x2014`.

A Unicode character can also be represented in a regular-expression by using its **Hex notation**(hexadecimal code point value) directly as described in construct `\x{...}`, for example a supplementary character U+2011F can be specified as `\x{2011F}`, instead of two consecutive Unicode escape sequences of the surrogate pair `\uD840\uDD1F`.

Unicode scripts, blocks, categories and binary properties 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.

Scripts, blocks, categories and binary properties can be used both inside and outside of a character class.

**Scripts** are specified either with the prefix `Is`, as in `IsHiragana`, or by using the `script` keyword (or its short form `sc`) as in `script=Hiragana` or `sc=Hiragana`.

The script names supported by `Pattern` are the valid script names accepted and defined by `UnicodeScript.forName`.

**Blocks** are specified with the prefix `In`, as in `InMongolian`, or by using the keyword `block` (or its short form `blk`) as in `block=Mongolian` or `blk=Mongolian`.

The block names supported by `Pattern` are the valid block names accepted and defined by `UnicodeBlock.forName`.

**Categories** may be specified with the optional prefix `Is`: Both `\p{L}` and `\p{IsL}` denote the category of Unicode letters. Same as scripts and blocks, categories can also be specified by using the keyword `general_category` (or its short form `gc`) as in `general_category=Lu` or `gc=Lu`.

The supported categories are those of [The Unicode Standard](#) in the version specified by the `Character` class. The category names are those defined in the Standard, both normative and informative.

**Binary properties** are specified with the prefix `Is`, as in `IsAlphabetic`. The supported binary properties by `Pattern` are

- Alphabetic
- Ideographic
- Letter
- Lowercase
- Uppercase
- Titlecase
- Punctuation
- Control
- White\_Space
- Digit
- Hex\_Digit
- Noncharacter\_Code\_Point
- Assigned

**Predefined Character classes** and **POSIX character classes** are in conformance with the recommendation of *Annex C: Compatibility Properties of [Unicode Regular Expression](#)*, when `UNICODE_CHARACTER_CLASS` flag is specified.

Classes	Matches
<code>\p{Lower}</code>	A lowercase character: <code>\p{IsLowercase}</code>
<code>\p{Upper}</code>	An uppercase character: <code>\p{IsUppercase}</code>
<code>\p{ASCII}</code>	All ASCII: <code>[\x00-\x7F]</code>
<code>\p{Alpha}</code>	An alphabetic character: <code>\p{IsAlphabetic}</code>
<code>\p{Digit}</code>	A decimal digit character: <code>\p{IsDigit}</code>
<code>\p{Alnum}</code>	An alphanumeric character: <code>[\p{IsAlphabetic}\p{IsDigit}]</code>
<code>\p{Punct}</code>	A punctuation character: <code>\p{IsPunctuation}</code>
<code>\p{Graph}</code>	A visible character: <code>[^\p{IsWhite_Space}\p{gc=Cc}\p{gc=Cs}\p{gc=Cn}]</code>
<code>\p{Print}</code>	A printable character: <code>[\p{Graph}\p{Blank}&amp;&amp;[^\p{Cntrl}]]</code>
<code>\p{Blank}</code>	A space or a tab: <code>[\p{IsWhite_Space}&amp;&amp;[^\p{gc=Zl}\p{gc=Zp}\x0a\x0b\x0c]</code>
<code>\p{Cntrl}</code>	A control character: <code>\p{gc=Cc}</code>
<code>\p{XDigit}</code>	A hexadecimal digit: <code>[\p{gc=Nd}\p{IsHex_Digit}]</code>
<code>\p{Space}</code>	A whitespace character: <code>\p{IsWhite_Space}</code>
<code>\d</code>	A digit: <code>\p{IsDigit}</code>
<code>\D</code>	A non-digit: <code>[^\d]</code>
<code>\s</code>	A whitespace character: <code>\p{IsWhite_Space}</code>
<code>\S</code>	A non-whitespace character: <code>[^\s]</code>
<code>\w</code>	A word character: <code>[\p{Alpha}\p{gc=Mn}\p{gc=Me}\p{gc=Mc}\p{Digit}\p{gc=Nd}\p{gc=No}]</code>

`\W` A non-word character: `[^\w]`

Categories that behave like the `java.lang.Character` boolean is *methodname* methods (except for the deprecated ones) are available through the same `\p{prop}` syntax where the specified property has the name `javamethodname`.

## Comparison to Perl 5

The `Pattern` engine performs traditional NFA-based matching with ordered alternation as occurs in Perl 5.

Perl constructs not supported by this class:

- Predefined character classes (Unicode character)

`\h` A horizontal whitespace

`\H` A non horizontal whitespace

`\v` A vertical whitespace

`\V` A non vertical whitespace

`\R` Any Unicode linebreak

sequence `\u000D\u000A|[\u000A\u000B\u000C\u000D\u0085\u2028\u2029]`

`\X` Match Unicode [extended grapheme cluster](#)

- The backreference constructs, `\g{n}` for the *n*<sup>th</sup> [capturing group](#) and `\g{name}` for [named-capturing group](#).
- The named character construct, `\N{name}` for a Unicode character by its name.
- The conditional constructs `(?(condition) X)` and `(?(condition) X| Y)`,
- The embedded code constructs `(?{code})` and `(??{code})`,
- The embedded comment syntax `(?#comment)`, and
- The preprocessing operations `\l`, `\u`, `\L`, and `\U`.

Constructs supported by this class but not by Perl:

- 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.



For a more precise description of the behavior of regular expression constructs, please see [\*Mastering Regular Expressions, 3rd Edition\*, Jeffrey E. F. Friedl, O'Reilly and Associates, 2006.](#)