## Summary of regular-expression constructs

| Construct | Matches |
|---|---|

### Characters

| Construct | Matches |
|---|---|
| *x* | The character *x* |
| \\ | The backslash character |
| \0*n* | The character with octal value 0*n* (0 <= *n* <= 7) |
| \0*nn* | The character with octal value 0*nn* (0 <= *n* <= 7) |
| \0*mnn* | The character with octal value 0*mnn* (0 <= *m* <= 3, 0 <= *n* <= 7) |
| \x*hh* | The character with hexadecimal value 0x*hh* |
| \u*hhhh* | The character with hexadecimal value 0x*hhhh* |
| \t | The tab character ('\u0009') |
| \n | The newline (line feed) character ('\u000A') |
| \r | The carriage-return character ('\u000D') |
| \f | The form-feed character ('\u000C') |
| \a | The alert (bell) character ('\u0007') |
| \e | The escape character ('\u001B') |
| \c*x* | The control character corresponding to *x* |

### Character classes

| Construct | Matches |
|---|---|
| [abc] | a, b, or c (simple class) |
| [^abc] | Any character except a, b, or c (negation) |
| [a-zA-Z] | a through z or A through Z, inclusive (range) |
| [a-d[m-p]] | a through d, or m through p: [a-dm-p] (union) |
| [a-z&&[def]] | d, e, or f (intersection) |
| [a-z&&[^bc]] | a through z, except for b and c: [ad-z] (subtraction) |
| [a-z&&[^m-p]] | a through z, and not m through p: [a-lq-z] (subtraction) |

### Predefined character classes

| Construct | Matches |
|---|---|
| . | Any character (may or may not match line terminators) |
| \d | A digit: [0-9] |
| \D | A non-digit: [^0-9] |
| \s | A whitespace character: [ \t\n\x0B\f\r] |
| \S | A non-whitespace character: [^\s] |
| \w | A word character: [a-zA-Z_0-9] |
| \W | A non-word character: [^\w] |

### POSIX character classes (US-ASCII only)

| Construct | Matches |
|---|---|
| \p{Lower} | A lower-case alphabetic character: [a-z] |
| \p{Upper} | An upper-case alphabetic character: [A-Z] |

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

## Classes for Unicode blocks and categories

| | |
|---|---|
| `\p{InGreek}` | A character in the Greek block (simple block) |
| `\p{Lu}` | An uppercase letter (simple category) |
| `\p{Sc}` | A currency symbol |
| `\P{InGreek}` | Any character except one in the Greek block (negation) |
| `[\p{L}&&[^\p{Lu}]]` | Any letter except an uppercase letter (subtraction) |

## Boundary matchers

| | |
|---|---|
| `^` | 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 terminator, 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 <u>capturing group</u> |

### Back references

| | |
|---|---|
| \\*n* | Whatever the *n*<sup>th</sup> <u>capturing group</u> 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 (non-capturing)

| | |
|---|---|
| (?:*X*) | *X*, as a non-capturing group |
| (?idmsux-idmsux) | Nothing, but turns match flags on - off |
| (?idmsux-idmsux:*X*) | *X*, as a <u>non-capturing group</u> with the given flags 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 or other character escapes. 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:

| | | |
|---|---|---|
| **1** | Literal escape | `\x` |
| **2** | Grouping | `[...]` |
| **3** | Range | `a-z` |
| **4** | Union | `[a-e][i-u]` |
| **5** | Intersection | `[a-z&&[aeiou]]` |

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

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 `"b"`. 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.