

APPENDIX A

Regular Expression Reference

THIS APPENDIX PROVIDES a comprehensive quick reference for your day-to-day Java regular expression needs. The material in this appendix is presented as pure regex patterns, not the Java `String`-delimited counterparts. For example, when referring to a digit, `\d`, not `\\d`, is used here.

Table A-1. Common Characters

| Regex | Description | Notes |
|-----------------|--|---|
| <code>q</code> | The character <i>q</i> | Could be any character, not including special regex characters or punctuation. |
| <code>\\</code> | The backslash (<code>/</code>) character | <code>\</code> delimits special regex characters, of which the backslash is a member. |
| <code>\t</code> | The tab character | |
| <code>\n</code> | The newline or linefeed character | |
| <code>\r</code> | The carriage-return character | |
| <code>\f</code> | The form-feed character | |

Table A-2. Predefined Character Classes

| Regex | Description | Notes |
|-----------------|--|---|
| | Any single character | Matches any single character, including spaces, punctuation, letters, or numbers. It may or may not match line-termination characters, depending on the operating system involved, or if the DOTALL flag is active in the regex pattern. Thus, it's probably a good idea to explicitly set, or turn off, DOTALL support in your patterns, in case you need to port your code. |
| <code>\d</code> | Any single digit from 0 to 9 | |
| <code>\D</code> | Will match any character except a single digit | By default, this won't match line terminators. |
| <code>\s</code> | A whitespace character: <code>[\t\n\x0B\f\r]</code> | This matches tab, space, end-of-line, carriage-return, form-feed, and newline characters. |
| <code>\S</code> | A non-whitespace character | Matches anything that is not a whitespace character, as described previously. Thus, <code>7</code> would be matched, as would punctuation. |
| <code>\w</code> | A word character: <code>[a-zA-Z_0-9]</code> | Any uppercase or lowercase letter, digit, or the underscore character. |
| <code>\W</code> | A nonword character; the opposite of <code>\w</code> | Anything that isn't a word character, as described previously. Thus, the minus sign will match, as will a space. It won't match the end-of-line <code>\$</code> or beginning-of-line <code>^</code> characters. |

Table A-3. Character Classes

| Regex | Description | Notes |
|------------------------------------|---|---|
| <code>[abc]</code> | <i>a</i> , <i>b</i> , or <i>c</i> | Strictly speaking, it won't match <i>ab</i> . |
| <code>[^abc]</code> | Any character except <i>a</i> , <i>b</i> , or <i>c</i> | This will match any character except <i>a</i> , <i>b</i> , or <i>c</i> . It won't match the end-of-line <code>\$</code> or beginning-of-line <code>^</code> characters. |
| <code>[a-zA-Z]</code> | Any uppercase or lowercase letter | When working with numbers, <code>[0-25]</code> doesn't mean 0 to 25. It means 0 to 2, or just 5. If you wanted 0 to 25, you would need to actually write an expression, such as <code>\d 1\d 2[0-5]</code> . Note that <code>[0-9]</code> is exactly equal to <code>\d</code> . |
| <code>[a-c[x-z]]</code> | <i>a</i> through <i>c</i> , or <i>x</i> through <i>z</i> | For example, <code>[1-3[7-9]]</code> matches 1 through 3, or 7 through 9. No other digit will do. |
| <code>[a-z&&[a,e]]</code> | <i>a</i> or <i>e</i> | <code>[a-z&&[a,e,i,o,u]]</code> matches all lowercase vowels. |
| <code>[a-z&&[^bc]]</code> | All lowercase letters except for <i>b</i> and <i>c</i> | For example, all the prime numbers between 1 and 9 would be <code>[1-9&&[^4689]]</code> . That is, 1 through 9, excluding 4, 6, 8, and 9. |
| <code>[a-d&&[^b-c]]</code> | <i>a</i> through <i>d</i> , but not <i>b</i> through <i>c</i> | <code>[1-9&&[^4-6]]</code> matches 1 through 3, or 7 through 9. Compare this to the union example presented earlier in this table. |

Table A-4. POSIX Character Classes

| Regex | Description | Notes |
|-------------------------|-----------------------------------|--|
| <code>\p{Lower}</code> | A lowercase alphabetic character | |
| <code>\p{Upper}</code> | An uppercase alphabetic character | |
| <code>\p{ASCII}</code> | An ASCII character | |
| <code>\p{Alpha}</code> | An alphabetic character | |
| <code>\p{Digit}</code> | A decimal digit | |
| <code>\p{Alnum}</code> | An alphanumeric character | |
| <code>\p{Punct}</code> | Punctuation | This is a good way to deal with punctuation in general, without having to delimit special characters such as periods, parentheses, brackets, and such. It matches <code>!"#\$%&'()*+,-./:;<=>?@[\\]^_`{ }~.</code> |
| <code>\p{Graph}</code> | A visible character | Exactly equal to <code>[\p{Alnum}\p{Punct}]</code> . |
| <code>\p{Print}</code> | A printable character | Exactly equal to <code>[\p{Graph}]</code> . |
| <code>\p{Blank}</code> | A space or a tab | |
| <code>\p{Space}</code> | Any whitespace character | It matches <code>[\t\n\x0B\f\r]</code> . |
| <code>\p{Cntrl}</code> | A control character | |
| <code>\p{XDigit}</code> | A hexadecimal digit | |

Table A-5. Boundary Matchers

| Regex | Description | Notes |
|-----------|--|--|
| ^ | Beginning-of-line character | This is an invisible character. |
| \$ | End-of-line character | This is an invisible character. |
| \b | A word boundary | <p>This is the position of a word boundary. Its usage requires some caution, because \b doesn't match characters; it matches a position. Thus, the String <i>anna marie</i> doesn't match the regex <i>anna\bmarie</i>. However, it does match <i>anna\b\smarrie</i>. That's because there's a word boundary at the position after the last <i>a</i> in <i>anna</i>, and it happens to be the space character, so \s is necessary to match it, and <i>marie</i> must then follow it. Because \b matches a position, it is meaningless to add greedy qualifiers to it. Thus, \b+, \b\b\b\b\b\b, and \b all match exactly the same thing.</p> <p>Further complicating the picture is the fact that in a character class, \b means a backspace character. This is syntactically legal (if a little awkward) because the word boundary has no place inside a character class. Thus, [\b] describes a backspace character, because it is surrounded by [and].</p> |
| \B | A non-word boundary | This is the opposite of word boundary, as described previously. |
| \A | The beginning of the input | \A matches the beginning of the input, but it isn't just a synonym for the ^ pattern. This distinction becomes clear if you use the <code>Pattern.MULTILINE</code> flag when you compile your pattern. \A matches the beginning of the input, which is the very beginning of the file. By contrast, ^ matches the beginning of each line when the <code>Pattern.MULTILINE</code> flag is active. |
| \Z | The end of the input except for the final \$, if any | \Z matches the end of the input, but it isn't just a synonym for the \$ character. This distinction becomes clear if you use the <code>Pattern.MULTILINE</code> flag when you compile your pattern. \Z matches the end of the input, which is the very end of the file. By contrast, \$ matches the end of each line when the <code>Pattern.MULTILINE</code> flag is active. |
| \G | The end of the previous match | |
| \z | The end of the input | This behaves exactly like the \Z with a capital <i>Z</i> character, except that it also captures the closing \$ character. |

Table A-6. Greedy Quantifiers

| Regex | Description | Notes |
|---------------|---|---|
| X? | X, once or not at all | A? would match <i>A</i> , or the absence of <i>A</i> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X* | X, zero or more times | This pattern is very much like the ? pattern, except that it matches zero or more occurrences. It doesn't match "any character," as its usage in DOS might indicate. Thus, A* would match <i>A</i> , <i>AA</i> , <i>AAA</i> , or the absence of <i>A</i> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X+ | X, one or more times | This quantifier is very much like the * pattern, except that it looks for the existence of one or more occurrences instead of zero or more occurrences. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X{n} | X, exactly <i>n</i> times | This quantifier demands the occurrence of the target exactly <i>n</i> times. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X{n,} | X, at least <i>n</i> times | This quantifier demands the occurrence of the target at least <i>n</i> times. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X{n,m} | X, at least <i>n</i> but not more than <i>m</i> times | This quantifier demands the occurrence of the target at least <i>n</i> times, but not more than <i>m</i> times. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-7. Reluctant Quantifiers

| Regex | Description | Notes |
|------------|-----------------------|---|
| X?? | X, once or not at all | This pattern is very much like the <code>?</code> pattern, except that it prefers to match nothing at all. When it's used with the <code>Matcher.matches()</code> method, <code>??</code> functions in exactly the same way as the <code>?</code> pattern. However, when it's used with <code>Matcher.find()</code> , the behavior is different. For example, the pattern <code>x??</code> , as applied to the String <code>xx</code> , will actually not find <code>x</code> , yet consider that lack of finding a success. That's because we asked it to be reluctant to match, and the most reluctant thing it can do is match zero occurrences of <code>x</code> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X*? | X, zero or more times | This pattern is very much like the <code>*</code> pattern, except that it prefers to match as little as possible. When it's used with the <code>Matcher.matches()</code> method, <code>*?</code> functions in exactly the same way as the <code>*</code> pattern. However, when it's used with <code>Matcher.find()</code> , the behavior is different. For example, the pattern <code>x*?</code> , as applied to the String <code>xx</code> , will actually not find <code>x</code> , yet consider that a success. That's because we asked it to be reluctant to match, and the most reluctant thing it can do is match zero occurrences of <code>x</code> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if character the class immediately precedes it). |
| X+? | X, one or more times | This pattern is very much like the <code>+</code> pattern, except that it prefers to match as little as possible. When it's used with the <code>Matcher.matches()</code> method, <code>+?</code> functions in exactly the same way as the <code>+</code> pattern. However, when it's used with <code>Matcher.find()</code> , the behavior is different. For example, the pattern <code>x+?</code> , as applied to the String <code>xx</code> , will actually find one <code>x</code> , yet consider that a success. That's because we asked it to be reluctant to match, and the most reluctant thing it can do is match one occurrence of <code>x</code> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-7. Reluctant Quantifiers (Continued)

| Regex | Description | Notes |
|----------------------|---|--|
| <code>X{n}?</code> | <i>X</i> , exactly <i>n</i> times | This pattern is exactly like the <code>X{n}</code> pattern. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| <code>X(n,)?</code> | <i>X</i> , at least <i>n</i> times | This pattern is very much like the <code>X{n,}</code> pattern, except that it prefers to match as little as possible. When it's used with the <code>Matcher.matches()</code> method, <code>X(n,)?</code> functions in exactly the same way as the <code>X{n,}</code> pattern. However, when it's used with <code>Matcher.find()</code> , the behavior is different. For example, the pattern <code>X{3,}?</code> , as applied to the <code>String xxxxx</code> , will actually only find <code>xxx</code> , yet consider that a success. Compare this with just <code>X{3,5}</code> , which would have found <code>xxxxx</code> . That's because we asked it to be reluctant to match, and the most reluctant thing it can do is match three occurrences of <i>x</i> . This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| <code>X{n,m}?</code> | <i>X</i> , at least <i>n</i> but not more than <i>m</i> times | This pattern is very much like the <code>X{n,m}</code> pattern, except that it prefers to match as little as possible. When it's used with the <code>Matcher.matches()</code> method, <code>X{n,m}?</code> functions in exactly the same way as the <code>X{n,m}</code> pattern. However, when it's used with <code>Matcher.find()</code> , the behavior is distinct. For example, the pattern <code>X{3,5}?</code> , as applied to the <code>String xxxxx</code> , will actually find <code>xxx</code> , yet consider that a success. Compare this with just <code>X{3,5}</code> , which would have found <code>xxxxx</code> . This happens because we asked it to be reluctant to match, and the most reluctant thing it can do is match three occurrences of <i>x</i> . Notice that if there were six <i>x</i> characters, such as <code>xxxxxx</code> , the pattern would have matched twice: once for the first three <code>xxx</code> characters and then again for the other three. Again, this is because three is the minimum requirement. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-8. Possessive Quantifiers

| Regex | Description | Notes |
|------------|-----------------------|---|
| X?+ | X, once or not at all | Very much like the <code>?</code> pattern, this pattern prefers to match as much as possible. However, this pattern won't release matches to help the entire expression match as a whole. For example, the pattern <code>\w?+\d</code> , as applied to the String <code>A2</code> , will actually not match, because the first <code>\w?+</code> consumes the <code>A</code> and the <code>2</code> , and it won't release them for the greater good of allowing the entire expression to match. Thus, <code>\d</code> is unable to match. This is because we asked <code>\w?+</code> to be possessive, and the most possessive thing it can do is match the occurrence of <code>A</code> and <code>2</code> , and not release them. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X*+ | X, zero or more times | Very much like the <code>*</code> pattern, this pattern prefers to match as much as possible. However, this pattern won't release matching to help the entire expression as a whole match. For example, the pattern <code>\w*+\d</code> , as applied to the String <code>Java2</code> , will actually not match, because the first <code>\w*+</code> consumes the String <code>Java2</code> and won't release it for the greater good of allowing the entire expression to match. Thus, <code>\d</code> is unable to match. This is because the pattern <code>\w*+</code> is possessive, and the most possessive thing it can do is match the entire <code>Java2</code> and not release anything. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-8. Possessive Quantifiers (Continued)

| Regex | Description | Notes |
|---------------|------------------------------------|---|
| X++ | <i>X</i> , one or more times | Very much like the <code>+</code> pattern, this pattern prefers to match as much as possible. However, this pattern won't release matching to help the entire expression as a whole match. For example, the pattern <code>\w++\d</code> , as applied to the String <i>Java2</i> , will actually not match, because the first <code>\w++</code> consumes the String <i>Java2</i> and won't release it for the greater good of allowing the entire expression to match. Thus, <code>\d</code> is unable to match. This is because the pattern <code>\w++</code> is possessive, and the most possessive thing it can do is match the entire <i>Java2</i> and not release anything. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X{n}+ | <i>X</i> , exactly <i>n</i> times | This pattern is exactly like the <code>X{n}</code> pattern. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |
| X{n,}+ | <i>X</i> , at least <i>n</i> times | Very much like the <code>X{n,}</code> pattern, this pattern prefers to match as much as possible. However, this pattern won't release matching to help the entire expression as a whole match. For example, the pattern <code>\w{4,}+\d</code> , as applied to the String <i>Java2</i> , will actually not match, because the <code>\w{4,}+</code> consumes the String <i>Java2</i> and won't release it for the greater good of allowing the entire expression to match. Thus, <code>\d</code> is unable to match 4. This is because the pattern <code>\w{4,}+</code> is possessive, and the most possessive thing it can do is match the entire <i>Java2</i> and not release anything. This pattern applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-8. Possessive Quantifiers (Continued)

| Regex | Description | Notes |
|----------------|---|---|
| X{n,m}+ | X, at least <i>n</i> but not more than <i>m</i> times | Very much like the <i>X{n,m}</i> pattern, this pattern prefers to match as much as possible. However, this pattern won't release matching to help the entire expression as a whole match. For example, the pattern <i>\x{2,5}\d</i> , as applied to the String <i>Java2</i> , will actually not match, because the first <i>\w++</i> consumes the String <i>Java2</i> , and won't release it for the greater good of allowing the entire expression to match. Thus, <i>\d</i> is unable to match. This is because the pattern <i>\x{2,5}</i> is possessive, and the most possessive thing it can do is match the entire <i>Java2</i> and not release anything. This applies to either the character that immediately precedes it, or a group (if the group immediately precedes it), or a character class (if the character class immediately precedes it). |

Table A-9. Logical Operators

| Regex | Description | Notes |
|------------|-----------------|---|
| XY | X followed by Y | This is the default relationship assumption between characters. Note that spaces are a valid part of this syntax. Thus, <i>A B</i> means the character <i>A</i> , followed by a space, followed by the character <i>B</i> . |
| X Y | Either X or Y | <i>AB CD</i> will match either <i>AB</i> or <i>CD</i> . Similarly, the pattern <i>hello sir madam</i> will match <i>hello sir</i> or it will match <i>madam</i> . Specifically, it won't match <i>hello madam</i> . This is because of the nature of the And pattern discussed previously. When the regex engine sees <i>hello sir</i> , it assumes you mean that <i>hello</i> , followed by a space, followed by <i>sir</i> should be treated as a single logical unit. These are all Anded together. Then the engine sees the Or pattern, so it assumes that the logical alternative is <i>madam</i> . If you actually want to accept <i>hello sir</i> or <i>hello madam</i> , you'll have to use groups—thus, the pattern <i>hello (sir madam)</i> . Or better yet, you can use the noncapturing group <i>hello (?:sir madam)</i> . |

Table A-9. Logical Operators (Continued)

| Regex | Description | Notes |
|-----------------|---|---|
| (X) | X, the capturing group | A capturing group is a logical unit that is conceptually similar to the logical units you're familiar with from algebra. Thus, <code>(\d\d\d\d)</code> is a capturing group that defines four digits. Capturing groups can be referred to later in your expression by using a back reference, as explained later. They're counted left to right and can be nested. Thus, <code>h(ello (world))</code> has three capturing groups. Capturing group 0 is the entire expression, which matches the String <code>hello world</code> . Capturing group 1 is <code>ello world</code> , because you count from left to right, and the first group starts with the <code>(</code> right before the <code>e</code> in <code>hello</code> . Group 2 is <code>world</code> , because the second group starts right before the <code>w</code> in <code>world</code> . |
| <code>\n</code> | The <i>n</i> th capturing group matched | <p>In this context, I'm not referring to newline, even though <code>\n</code> looks like the newline symbol. The <i>n</i> in this case refers to a number. The regex engine allows you to access the information captured by a previous part of the group, even as the search is executing. For example, if you want to find repeated words, all you need is the pattern <code>(\w+)\W\1</code>, which says, "Look for a group of word characters, followed by a nonword character, followed by that exact same word character captured in group 1." If you attempt to refer to a group that doesn't exist, a <code>PatternSyntaxException</code> will be thrown.</p> <p>If you happen to have, say, 13 captured groups, then <code>\13</code> will mean that you want the thirteenth capturing group. If you don't have 13 groups, then the same expression <code>\13</code> will mean the first capturing group, followed by the digit 3.</p> |

Table A-10. Quotation

| Regex | Description | Notes |
|-----------------|---|--|
| <code>\</code> | Quotes the following character | This quotes the metacharacter that follows, so it will actually be treated as a character. Thus, if you were looking for a dollar sign, you would use <code>\\$</code> . as the pattern. By contrast, <code>\$</code> would have matched the end-of-line character. Remember that for regex expressions used directly as Strings, you need to double the number of <code>\</code> characters you see. Thus, in a Java String, <code>\s</code> becomes <code>\\s</code> . |
| <code>\Q</code> | Quotes all characters until <code>\E</code> | This works in conjunction with <code>\E</code> to quote a sequence of characters. If you need to quote a lot of characters in sequence, then use <code>\Q</code> to open your quote and <code>\E</code> to close it. For example, if you want the characters <code>\(/?*</code> , the expression <code>\Q\(/?*\\E</code> will do the job. |
| <code>\E</code> | Ends quote started by <code>\Q</code> | |

Table A-11. Noncapturing Group Constructs

| Regex | Description | Notes |
|---------------------------------|---|--|
| <code>(?:X)</code> | Defines a subpattern as a logical unit | Noncapturing groups don't store the information that actually matches the pattern for later access. These are much more efficient than capturing groups if you're only using grouping for logical purposes. This pattern is noncapturing. |
| <code>(?idmsux-idmsux)</code> | <i>i</i> for CASE_INSENSITIVE <i>x</i> for COMMENTS <i>s</i> for DOTALL <i>u</i> for UNICODE_CASE <i>m</i> for MULTILINE <i>d</i> for UNIX_LINES | The pattern <code>(?i)hel(?-i)LO</code> will match the String <code>HELLO</code> , because <code>(?i)</code> indicates a case-insensitive match starting from <i>h</i> , and <code>(?-i)</code> signals an end to that case insensitivity after the first <i>l</i> . This pattern is noncapturing. |
| <code>(?idmsux-idmsux:X)</code> | X, with the given flags on or off | The pattern <code>(?i:hel)LO</code> will match the String <code>HELLO</code> , because <code>(?i:</code> indicates a case-insensitive match starting from <i>h</i> and ending with the first <i>l</i> . This pattern is noncapturing. |

Table A-12. Lookarounds

| Regex | Description | Notes |
|--------|---|--|
| (?=X) | X, using zero-width positive lookahead | This pattern glances to the right of whatever remains to be parsed from the candidate String to find the first position at which the expression X exists. For example, if you want to extract all of the inline comments from a text file, you might try the pattern (?!//).*\$ and extract group 0. This pattern is noncapturing. |
| (?!X) | X, using zero-width negative lookahead | This pattern glances to the right, to whatever remains to be parsed from the candidate String, to find the first position at which the expression X doesn't exist. For example, if you want to skip leading spaces leading up to some content, you could use (?!\s).*. This pattern is noncapturing. |
| (?<=X) | X, using zero-width positive lookbehind | This pattern glances to the left, to whatever remains to be parsed from the candidate String, to find the first position at which the expression X exists. For example, if you want to extract all of the inline comments from a text file, you might try the pattern (?=//).*\$\$. This pattern is noncapturing. |
| (?<!X) | X, using zero-width negative lookbehind | This pattern glances to the left, to whatever remains to be parsed from the candidate String, to find the first position at which the expression X doesn't occur. For example, if you want to extract all of the text before inline Java comments from a text file, you might try the pattern .*(?<=//). This pattern is noncapturing. |
| (?>X) | X, as an independent, non-capturing group | This pattern refuses to release the contents of the match, regardless of the consequences on the rest of the pattern's ability to match. Thus, whereas the pattern <code>\w+\d</code> matches the String <code>java2</code> , the pattern <code>(?>\w+)\d</code> does not, because the <code>(?>\w+)</code> consumes <code>java</code> and <code>2</code> , and refuses to release the <code>2</code> so that <code>\d</code> can match. |

Table A-13. Less Common Characters

| Regex | Description | Notes |
|---------------------|--|--|
| <code>\0n</code> | The character with octal value <i>on</i> | $0 \leq n \leq 7$ |
| <code>\0nn</code> | The character with octal value <i>onn</i> | $0 \leq n \leq 7$ |
| <code>\0mnn</code> | The character with octal value <i>omnn</i> | $0 \leq m \leq 3, 0 \leq n \leq 7$ (This can't exceed 377.) |
| <code>\xhh</code> | The character with hexadecimal value <i>0xhh</i> | $0 \leq h \leq 9$ or $A \leq h \leq F$ |
| <code>\uhhhh</code> | The character with hexadecimal value <i>0xhhhh</i> | $0 \leq h \leq 9$ or $A \leq h \leq F$ |
| <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> | |

Table A-14. Unicode Blocks and Categories

| Regex | Description | Notes |
|---|--|--|
| <code>\p{InGreek}</code> | A character in the Greek block | |
| <code>\p{Lu}</code> | An uppercase letter | <code>\p{Lu}</code> matches any uppercase character. |
| <code>\p{Sc}</code> | A currency symbol | If you need to find or swap out, say, a dollar sign, this is a good way to do so without having to deal with the various delimiting complexities of not matching the end-of-line character \$. |
| <code>\P{InGreek}</code> | Any character except one in the Greek block (negation) | Notice the use of the capital <i>P</i> here. In general, uppercase <code>\P</code> is the opposite of lowercase <code>\p</code> . Thus, <code>\P{Lower}</code> matches all uppercase characters. |
| <code>[\p{L} & & [^ \p{Lu}]]</code> | Any non-uppercase letter | This is exactly equal to <code>\p{Upper}</code> . |

APPENDIX B

Pattern and Matcher Methods

THIS APPENDIX PROVIDES a summary of the methods of the `Pattern` and `Matcher` classes in Java. It's intended to be a quick reference for working with the various regex utilities you'll be using. For more detailed descriptions, please see the appropriate section in the text.

Pattern Class Fields

UNIX_LINES

The `UNIX_LINES` flag is used in constructing the second parameter of the `Pattern.compile(String regex, int flags)` method. Use this flag when parsing data that originates on a UNIX machine.

On many flavors of UNIX, the invisible character `\n` is used to note termination of a line. This is distinct from other operating systems, including flavors of Windows, which may use `\r\n`, `\n`, `\r`, `\u2028`, or `\u0085` for a line terminator.

If you've ever transported a file that originated on a UNIX machine to a Windows platform and opened it, you may have noticed that the lines sometimes don't terminate as you might expect, depending on which editor you use to view the text. This happens because the two systems can use different syntax to denote the end of the line.

The `UNIX_LINES` flag simply tells the regex engine that it's dealing with UNIX style lines, which affects the matching behavior of the regular expression meta-characters `^` and `$`.

NOTE Using the `UNIX_LINES` flag, or the equivalent `(?d)` regex pattern, doesn't degrade performance. By default, this flag isn't set.

CASE_INSENSITIVE

The `CASE_INSENSITIVE` field is used in constructing the second parameter of the `Pattern.compile(String regex, int flags)` method. It's useful when you need to match U.S. ASCII characters, regardless of case.

NOTE Using this flag, or the equivalent `(?i)` regular expression, can cause performance to degrade slightly. By default, this flag is not set.

COMMENTS

The `COMMENTS` field is defined because it's used in constructing the second parameter of the `Pattern.compile(String regex, int flags)` method. It tells the regex engine that the regex pattern has an embedded comment in it. Specially, it tells the regex engine to ignore any comments in the pattern, starting with the spaces leading up to the `#` character and everything thereafter, until the end of the line.

Thus, the regex pattern `A #matches uppercase US-ASCII char code 65` will use `A` as the regular expression, but the spaces leading up to the `#` character and everything thereafter until the end of the line will be ignored.

NOTE Using this flag, or the equivalent `(?x)` regular expression, doesn't degrade performance.

MULTILINE

The `MULTILINE` field is used in constructing the second parameter of the `Pattern.compile(String regex, int flags)` method. It tells the regex engine that regex input isn't a single line of code; rather, it contains several lines that have their own termination characters.

This means that the beginning-of-line character, `^`, and the end-of-line character, `$`, will potentially match several lines within the input `String`.

For example, imagine that your input String is *This is sentence.\n So is this*. If you use the MULTILINE flag to compile the regular expression pattern:

```
Pattern p = Pattern.compile("^", Pattern.MULTILINE);
```

then the beginning of line character, `^`, will match before the *T* in *This is a sentence*. It will also match just before the *S* in *So is this*. Without using the MULTILINE flag, the match will only find the *T* in *This is a sentence*.

NOTE Using this flag, or the equivalent *(?m)* regular expression, may degrade performance.

DOTALL

The DOTALL flag is used in constructing the second parameter of the `Pattern.compile(String regex, int flags)` method.

The DOTALL flag tells the regex engine to allow the metacharacter period (`.`) to match any character, *including* a line termination character. What does this mean?

Imagine that your candidate String were *Test\n*. If your corresponding regex pattern were the period (`.`), then you would normally have four matches: one for the *T*, another for the *e*, another for *s*, and the fourth for *t*. This is because the regex metacharacter period (`.`) will normally match any character, except line termination characters.

Enabling the DOTALL flag

```
Pattern p = Pattern.compile(".", Pattern.DOTALL);
```

would have generated five matches. Your pattern would have matched the *T*, *e*, *s*, and *t* characters. In addition, it would have matched the `\n` character at the end of the line.

NOTE Using this flag, or the equivalent *(?s)* regular expression, doesn't degrade performance.

UNICODE_CASE

The `UNICODE_CASE` flag in conjunction with the `CASE_INSENSITIVE` flag generates case-insensitive matches for international character sets.

NOTE Using this flag, or the equivalent *(?u)* regular expression, can degrade performance.

CANON_EQ

As you know, characters are actually stored as numbers. For example, in the U.S. ASCII character set, the character A is represented by the number 65. Depending on the character set that you're using, the same character can be represented by different numeric combinations. For example, à can be represented by both +00E0 and U+0061U+0300. A `CANON_EQ` match would match either representation.

NOTE Using this flag may degrade performance.

Pattern Class Methods

public static Pattern compile(String regex) throws PatternSyntaxException

You'll notice that the `Pattern` class doesn't have a public constructor. This means that you can't write the following type of code:

```
Pattern p = new Pattern("my regex");//wrong!
```

To get a reference to a `Pattern` object, you must use the static method `pattern(String regex)`. Thus, your first line of regex code might look like the following:

```
Pattern p = Pattern.compile("my regex");//Right!
```

The parameter for this method is a `String` that represents a regular expression. When passing a `String` to a method that expects a regular expression, it's important to delimit any `\` characters that the regular expressions might have by appending another `\` character to them. This is because `String` objects internally use the `\` character to delimit metacharacters in a character sequences, regardless of whether those character sequences are regular expressions. This has been true long before regular expression were a part of Java. Thus, the regular expression `\d` becomes `\\d`. To match a single digit, your regular expression code becomes the following:

```
Pattern p = Pattern.compile("\\d");
```

The point here is that the regular expression `\d` becomes the `String` `\\d`.

The delimitation of the `String` parameter can sometimes be tricky, so it's important to understand it well. By and large, it means that you double the `\` characters that might already be present in the regular expression. It doesn't mean that you simply append a single `\` character.

The `compile` method will throw a `java.util.regex.PatternSyntaxException` if the regular expression itself is badly formed. For example, if you pass in a `String` that contains `/4`, the `compile` method will throw a `PatternSyntaxException` at runtime, because the syntax of the regular expression `/4` is illegal.

The `compile(String regex)` method returns a `Pattern` object.

***public static compile pattern(String regex, int flags)
throws PatternSyntaxException***

The `compile(String regex, int flags)` method is a more powerful form of the `compile` method. The first parameter for this method, `regex`, is a `String` that represents a regular expression, as detailed in the previous `pattern.compile(String regex)` method entry. For details on how the `String` parameter must be formatted, please see the previous `compile(String regex)` method entry.

The flexibility of this `compile` method is fully realized by using the second parameter, `int flags`. For example, if you want a match to be successful regardless of the case of the candidate `String`, then your pattern might look like the following:

```
Pattern p = Pattern.compile(regex, Pattern.CASE_INSENSITIVE);
```

You can combine the flags by using the `|` operator. For example, to achieve case-insensitive Unicode matches that include a comment, you might use the following:

```
Pattern p =  
Pattern.compile("t # a compound flag example",Pattern.CASE_INSENSITIVE |  
                Pattern.UNICODE_CASE| Pattern.COMMENT);
```

The `compile(String regex, int flags)` method returns a `Pattern` object.

public String pattern()

This method returns the regular expression from which this pattern was compiled. This is a simple `String` that represents the regex you passed in.

This method can be misleading in two ways. First, the `String` that is returned doesn't reflect any flags that were set when the pattern was compiled. Second, the regex `String` you passed in isn't always the pattern `String` you get back out. Specifically, the original `String` delimitations aren't shown. Thus, if your original code was

```
Pattern p = Pattern.compile("\\d");
```

then you should expect your output to be `\d`, with a single `\` character.

public Matcher matcher(CharSequence input)

Remember that you create a `Pattern` object by compiling a description of what you're looking for. A `Pattern` lists the features of what you're looking for. Speaking purely conceptually, your patterns might look like the following:

```
Pattern p = Pattern.compile("She must have red hair and be smarter than I am");
```

Correspondingly, you'll need to compare that description against candidates. That is, you'll want to examine a given `String` to see whether it matches the description you provided.

The `Matcher` object is designed specifically to help you do this sort of interrogation. The `Pattern.matcher(CharSequence input)` method returns the `Matcher` that will help get details about how your candidate `String` compares with the description you passed in.

public int flags()

Earlier, I discussed the constant flags that you can use in compiling your regex pattern. The `flags` method simply returns an `int` that represents those flags.

To see if your `Pattern` class is currently using a given flag—for example, the `Pattern.COMMENTS` flag—simply extract the flag

```
int flgs = myPattern.flags();
```

and then & that flag to the `Pattern.COMMENTS` flag:

```
boolean isUsingCommentFlag = ( Pattern.COMMENTS == (Pattern.COMMENTS & flgs)) ;
```

Similarly, to see if you're using `CASE_insensitive`, do the following:

```
boolean isUsingCaseInsensitiveFlag =  
(Pattern.CASE_insensitive == (Pattern.CASE_insensitive & flgs));
```

public static boolean matches(String regex, CharSequence input)

Very often, you'll find that all you need to know about a `String` is whether it matches a given regular expression exactly. You don't want to have to create a `Pattern` object, extract its `Matcher` object, and interrogate that `Matcher`.

This static utility method is designed to help you do exactly that. Internally, it creates the `Pattern` and `Matcher` objects you need, compares the regex to the input `String`, and returns a `Boolean` indicating whether the two objects match exactly. Usage might look something like `PatternMatchesTest` example shown here:

```
import java.util.regex.*;  
public class PatternMatchesTest{  
    public static void main(String args[]){  
  
        String regex = "ad*";  
        String input = "add";  
  
        boolean isMatch = Pattern.matches(regex,input);  
        System.out.println(isMatch);\return true  
    }  
}
```

If you're going to be doing a lot of comparisons, then it's more efficient to explicitly create a `Pattern` object and do your matches manually. However, if you're not going to be doing a lot of comparisons, then `matches` is a handy utility method.

The `Pattern.matches(String regex, CharSequence input)` method is also used internally by the `String` class. As of J2SE 1.4, `String` has a new method called `matches`, which internally defers to this one. Thus, you might already be using this method without being aware of it.

Of course, this method can throw a `PatternSyntaxException` if the regex pattern under consideration isn't well formed.

`public String[] split(CharSequence input)`

This method can be particularly helpful if you need to break up a `String` into an array of substrings based on some criteria. In concept, it's similar to the `StringTokenizer`. However, it's much more powerful and more resource intensive than `StringTokenizer`, because it allows your program to use a regular expressions as the splitting criteria.

This method always returns at least one element. If the split candidate, input, can't be found, a `String` array is returned that contains exactly one `String`, namely the original input.

If the input can be found, then a `String` array is returned. That array contains every substring after an occurrence of the input. Thus, for the pattern

```
Pattern p = new Pattern.compile(",");
```

the `split` method for *Hello, Dolly* will return a `String` array consisting of two elements. The first element of the array will contain the `String` *Hello*, and the second will contain the `String` *Dolly*. That `String` array is obtained as follows:

```
String tmp[] = p.split("Hello,Dolly");
```

In this case, the value return is

```
//tmp =={ "Hello", "Dolly"}
```

There are some subtleties you should be aware of when working with this method. If the candidate `String` had been *Hello,Dolly*, with a trailing comma character after the *y* in *Dolly*, then this method would still have returned two elements: a `String` array consisting of *Hello* and *Dolly*. The implicit behavior is that training spaces aren't returned.

If the input `String` had been *Hello,,,Dolly*, the resulting `String` array would have four elements. The return value of the `split` method, as applied to the `Pattern`, is

```
// p.split("Hello,,,Dolly") == {"Hello", "", "", "Dolly"}
```

The `String` method further optimizes its search criteria by placing an invisible `^` before the pattern and a `$` after it.

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

This method works in exactly the same way as `Pattern.split(CharSequence input)`, with one variation. The second parameter, `limit`, allows you to control how many elements are returned:

```
Limit == 0
```

If you specify that the second parameter, `limit`, should equal 0, then this method behaves exactly like its overloaded counterpart:

```
Limit >0
```

Use a positive limit if you're interested in only a certain number of matches. You should use number *1* as the limit. Say the `Pattern p` has been compiled for the `String`, as previously. To split the `String` *Hello, Dolly, You, Are, My, Favorite* when you only want the first two tokens, you would use this:

```
String[] tmp = pattern.split("Hello, Dolly, You, Are, My, Favorite",3);
```

The value of the resulting `String` would be this:

```
//tmp[0] = "Hello", tmp[1] = "Dolly";
```

The interesting behavior here is that a third element is returned, in this case

```
//tmp[2] = "You, Are, My, Favorite";
```

Using a positive limit can potentially lead to performance enhancements, because the regex engine can stop searching when it meets the specified number of matches:

```
Limit <0
```

Using a negative number—any negative number—for the limit tells the regex engine that you want to return as many matches as possible *and* that you want trailing spaces, if any, to be returned. Thus, for the regex pattern `,` and the candidate `String` *Hello,Dolly* the command

```
String tmp[] = p.split("Hello,Dolly", -1);
```


results in

```
//tmp == {"Hello","Dolly"};
```

However, for the String *Hello, Dolly*,
spaces after the comma following the *Dolly*, the method call

```
String tmp[] = p.split("Hello,Dolly, ", -1);
```

results in

```
//tmp == {"Hello","Dolly"," "};
```

Notice that the actual value of the negative limit doesn't matter. Thus,

```
p.split("Hello,Dolly", -1);
```

is exactly equivalent to

```
p.split("Hello,Dolly", -100);
```

Method Class Methods

public Pattern pattern()

The `pattern` method returns the `Pattern` that created this particular `Matcher` object. The `Pattern` returned doesn't contain any of the flags that are explicitly set by using the `Pattern` constants when the pattern is compiled, such as `Pattern.MULTILINE`.

public Matcher reset()

The `reset()` method clears all state information from the `Matcher` object it is called on. The `Matcher` is, in effect, reverted to the state it originally had when you first received a reference to it.

public Matcher reset(CharSequence input)

The `reset(CharSequence input)` method clears the state of the `Matcher` object it's called on and replaces the candidate `String` with the new input. This has the same effect as creating a new `Matcher` object, except that it doesn't have the associated overhead. This recycling can be a useful optimization, and it's one that I often use.

public int start()

This method returns the index of the first character of the candidate `String` matched. If there are no matches, or if no matches have been attempted, an `IllegalStateException` is thrown.

public int start(int group)

This method allows you to specify which subgroup within a matching you're interested in and returns the index of the first character in which that subgroup starts. If there are no matches, or if no matches have been attempted, an `IllegalStateException` is thrown. If you refer to a group number that doesn't exist, an `IndexOutOfBoundsException` is thrown.

public int end()

The `end` method returns the ending index plus 1 of the last successful match the `Matcher` object had. If no matches exist, or if no matches have been attempted, this method throws an `IllegalStateException`.

public int end(int group)

Like the `start(int)` method, this method allows you to specify which subgroup within a matching you're interested in, except that it returns the last index of matching character sequence plus 1. If no matches exist, or if no matches have been attempted, this method throws an `IllegalStateException`. If you refer to a group number that doesn't exist, an `IndexOutOfBoundsException` is thrown.

public String group()

The `group` method can be a powerful and convenient tool in the war against jumbled code. It simply returns the substring of the candidate `String` that matches the original regex pattern. The `group()` method throws an `IllegalStateException` if the `find` method call wasn't successful. Similarly, it throws an `IllegalStateException` if `find` has never been called at all.

public String group(int group)

This method is a more powerful counterpart to the `group()` method. It allows you to extract parts of a candidate string that match a subgroup within your pattern.

The `group(int)` method throws an `IllegalStateException` if the `find` method call wasn't successful. Similarly, it throws an `IllegalStateException` if `find` has never been called at all. If called for a group number that doesn't exist, it throws an `IndexOutOfBoundsException`.

`public int groupCount()`

This method simply returns the number of groups that the Pattern defined. There's a very important, and somewhat counterintuitive, subtlety to notice about this method: It returns the number of possible groups based on the original Pattern, without even considering the candidate String. Thus, it's not really information about the Matcher object; instead, it's about the Pattern that helped spawn it. This can be tricky, because the fact that this method lives on the Matcher object could be interpreted as meaning that it's providing feedback about the state of the Matcher. It isn't. It's telling you how many matches are theoretically possible for the given Pattern.

`public boolean matches()`

This method is designed to help you match a candidate String against the matcher's Pattern. It returns `true` if, and only if, the candidate String under consideration matches the pattern exactly.

`public boolean find()`

The `find()` method parses just enough of the candidate String to find a match. If such a substring is found, then `true` is returned and `find` stops parsing the candidate. If no part of the candidate String matches the pattern, then `find` returns `false`.

`public boolean find(int start)`

`find(int start)` works exactly like its overloaded counterpart, with the exception of where it starts searching. The `int` parameter `start` simply tells the Matcher at which character to start its search.

Thus, for the candidate String *I love Java. Java is my favorite language. Java Java Java.* and the Pattern **Java**, if you only want to start searching at character index 11, you use the command `find(11)`.

```
public Matcher appendReplacement(StringBuffer sb,
String replacement)
```

The `appendReplace` method allows you to modify the contents of a `StringBuffer` based on a regular expression. It even allows you to use back references by using the `$n` notation, in which `n` refers to some captured subgroup. For example, for the `StringBuffer` *Waldo Smith*, the Pattern `(\w+) (\w+)`, and the replacement `$2, $1`, the contents of the `StringBuffer` will be modified to *Smith, Waldo*.

The `appendReplacement` method will throw an `IllegalStateException` if a `find()` hasn't been called, or if `find` would have returned `false` if called. It will throw an `IndexOutOfBoundsException` if the capturing group referred to by `$1`, `$2`, etc., doesn't exist in the part of the pattern currently being scrutinized by the `Matcher`.

```
public StringBuffer appendTail(StringBuffer sb)
```

The `appendTail` method is a supplement to the `appendReplacement` method. It simply appends every remaining subsequence from the original candidate `String` to the `StringBuffer`, if reading from the append position, which I explained in the `appendReplacement` section, to the end of the candidate string.

```
public String replaceAll(String replacement)
```

The `replaceAll` method returns a `String` that replaces every occurrence of the description with the replacement. Using this method will change the state of your `Matcher` object. Specifically, the `reset()` method will be called. Therefore, remember that all `start`, `end`, `group`, and `find` calls will have to be reexecuted.

Like the `appendReplacement` method, the `replaceAll` method can contain references to substring by using the `$` symbol. For details, please see the `appendReplacement` documentation presented earlier.

```
public String replaceFirst(String replacement)
```

The `replaceFirst` method is a more limited version of the `replaceAll` method. This method returns a `String` that replaces the first occurrence of the description with the replacement. Using this method will change the state of your `Matcher` object. Specifically, the `reset()` method will be called. Therefore, remember that all `start`, `end`, `group`, and `find` calls will have to be reexecuted after `replaceFirst` is called.

Like the `appendReplacement` method, the `replaceFirst` method can contain references to substring by using the `$` symbol. For details, please see the `appendReplacement` documentation presented earlier.

APPENDIX C

Common Regex Patterns

THIS APPENDIX PRESENTS some practical regex patterns that you can use for common matching and validation tasks.

Table C-1. IP Address `^(([0-1]?\d{1,2}\.)([2[0-4]\d\.)|(25[0-5]\.)){3}([0-1]?\d{1,2})|([2[0-4]\d)|([25[0-5]))$`

| Regex | Description |
|---------------------|---|
| <code>^</code> | Beginning of line |
| <code>(</code> | A group consisting of |
| <code>(</code> | A subgroup consisting of |
| <code>[0-1]?</code> | Zero or one, both optional, followed by |
| <code>\d</code> | Any digit |
| <code>{1,2}</code> | Repeated one or two times, followed by |
| <code>\.</code> | A period |
| <code>)</code> | Close subgroup |
| <code> </code> | Or |
| <code>(</code> | A subgroup consisting of |
| <code>2</code> | The digit 2, followed by |
| <code>[0-4]</code> | Any digit from 0 to 4, followed by |
| <code>\d</code> | Any digit, followed by |
| <code>\.</code> | A period |
| <code>)</code> | Close subgroup |
| <code> </code> | Or |
| <code>(</code> | A subgroup consisting of |
| <code>2</code> | The digit 2, followed by |
| <code>5</code> | The digit 5, followed by |

Table C-1. IP Address $\wedge(((0-1)?\backslash d\{1,2\}\backslash.)(2[0-4]\backslash d\backslash.)(25[0-5]\backslash.))\{3\}(((0-1)?\backslash d\{1,2\})|(2[0-4]\backslash d)|(25[0-5]))\$$ (Continued)

| Regex | Description |
|----------------|---|
| [0-5] | Any digit from 0 to 5, followed by |
| \. | A period |
|) | Close subgroup |
|) | Close group |
| {3} | Repeated three times, followed by |
| (| A group consisting of |
| (| A subgroup consisting of |
| [0-1]? | Zero or one, both optional, followed by |
| \d{1,2} | Any two digits |
|) | Close subgroup |
| | Or |
| (| A subgroup consisting of |
| 2 | The digit 2, followed by |
| [0-4] | Any digit from 0 to 4, followed by |
| \d | Any digit |
|) | Close subgroup |
| | Or |
| (| A subgroup consisting of |
| 2 | The digit 2, followed by |
| 5 | The digit 5, followed by |
| [0-5] | Any digit from 0 to 5, followed by |
|) | Close subgroup |
|) | Subgroup |
| \$ | End of line |

*** In English:** Three sets of three digits separated by periods, each ranging from 0 to 255, followed by three digits, each ranging from 0 to 255.

NOTE The following pattern also matches the IP address, but all the groups have been marked as noncapturing:

```
(?:([0-1]?[0-9]{1,2}\.)(?:2[0-4]\.)(?:25[0-5]\.)){3}(?:[0-1]?[0-9]{1,2})|(?:2[0-4]\.)(?:25[0-5]\.))
```

This is slightly more efficient than the previous pattern, but it's less legible.

Table C-2. Simple E-mail `^(\\p{Alnum}+\\.|_|\\-)?*\\p{Alnum}@\\p{Alnum}+\\.|_|\\-)?*\\p{Alpha}$`

| Regex | Description |
|-------------------------|--|
| <code>^</code> | Beginning of line |
| <code>(</code> | A group consisting of |
| <code>\\p{Alnum}</code> | A letter or a digit |
| <code>+</code> | Repeated one or more times, followed by |
| <code>(</code> | A subgroup consisting of |
| <code>\\.</code> | A period |
| <code> </code> | Or |
| <code>_</code> | An underscore |
| <code> </code> | Or |
| <code>\\-</code> | A hyphen |
| <code>)</code> | Close subgroup |
| <code>?</code> | The preceding punctuation is optional |
| <code>)</code> | Close group |
| <code>*</code> | The previous group can be repeated zero or more times, followed by |
| <code>\\p{Alnum}</code> | A letter or a digit, followed by |
| <code>@</code> | An @ symbol, followed by |
| <code>\\p{Alnum}</code> | A letter or a digit |
| <code>+</code> | Repeated one or more times, followed by |

Table C-2. Simple E-mail $^{\wedge}(\backslash p\{Alnum\}+(\backslash .|\backslash _|\backslash -)?)^*\backslash p\{Alnum\}@\backslash p\{Alnum\}+(\backslash .|\backslash _|\backslash -)?)^*\backslash p\{Alpha\}\$$ (Continued)

| Regex | Description |
|-----------|---------------------------------------|
| (| A group consisting of |
| \. | A period |
| | Or |
| _ | An underscore |
| | Or |
| \- | A hyphen |
|) | Close group |
| ? | The preceding punctuation is optional |
| \p{Alpha} | An upper- or lowercase letter |
| \$ | End of line |

* **In English:** Any number of alphanumeric characters followed by single hyphens, periods, or underscores, but ending in an alphanumeric character; followed by an @ symbol; followed by any number of alphanumeric characters; followed by single hyphens, periods, or underscores, but ending in an upper- or lowercase character.

NOTE The following pattern matches the previous one exactly, except that it allows an IP address as well:

```
^(\p{Alnum}+(\backslash .|\backslash _|\backslash -)?)^*\p{Alnum}@\(((\p{Alnum}+(\backslash .|\backslash _|\backslash -)?)^*\p{Alpha})|(((0-1)?\d{1,2}\.)(2[0-4]\d\.)(25[0-5]\.)){3}(((0-1)?\d{1,2})|(2[0-4]\d)|(25[0-5])))\)$
```

For a breakdown of the IP address pattern, please see Table C-1.

Table C-3. Digit Repeated Exactly n Times $\backslash d\{n\}$, Where n Is the Number of Digits Needed

| Regex | Description |
|---|--------------------|
| $\backslash d$ | Any number |
| $\{n\}$ | Repeated n times |
| * In English: n digits. Thus, if n was equal to 4, any four digits. | |

Table C-4. Characters Repeated Exactly n Times $\backslash w\{n\}$, Where n Is the Number of Characters Needed

| Regex | Description |
|---|--|
| $\backslash w$ | Any number, any digit, or an underscore symbol |
| $\{n\}$ | Repeated n times |
| * In English: n characters. Thus, if n was equal to 4, any four characters. | |

Table C-5. Characters Repeated n to m Times $\backslash w\{n..m\}$, Where n Is the Number of Characters Needed

| Regex | Description |
|--|--|
| $\backslash w$ | Any number, any digit, or an underscore symbol |
| $\{n$ | Repeated n times |
| $m\}$ | But not more than m times |
| * In English: n characters. Thus, if n was equal to 4 and m was equal to 9, any four, five, six, seven, eight, or nine characters. | |

Table C-6. Credit Cards: Visa, MasterCard, American Express, and Discover
`^((4\d{3})|(5(1-5)\d{2})|(6011))-?\d{4}-?\d{4}-?\d{4}|3[4,7]\d{13}$`

| Regex | Description |
|--------------------|--|
| <code>^</code> | Beginning of line |
| <code>(</code> | A group consisting of |
| <code>(</code> | A subgroup consisting of |
| <code>4</code> | the digit four, followed by |
| <code>\d{3}</code> | Any three digits |
| <code>)</code> | Close subgroup |
| <code> </code> | Or |
| <code>(</code> | A subgroup consisting of |
| <code>5</code> | The digit 5, followed by |
| <code>[1-5]</code> | Any digit ranging from 1 to 5, followed by |
| <code>\d{2}</code> | Any two digits |
| <code>)</code> | Close subgroup |
| <code> </code> | Or |
| <code>(</code> | A subgroup consisting of |
| <code>6001</code> | The digits 6, 0, 0, 1 |
| <code>)</code> | Close subgroup, followed by |
| <code>-?</code> | An optional hyphen, followed by |
| <code>\d{4}</code> | Any four digits, followed by |
| <code>-?</code> | An optional hyphen, followed by |
| <code>\d{4}</code> | Any four digits, followed by |
| <code>-?</code> | An optional hyphen, followed by |
| <code>\d{4}</code> | Any four digits, followed by |
| <code>-?</code> | An optional hyphen |
| <code> </code> | Or |
| <code>3</code> | The digit 3, followed by |

Table C-6. Credit Cards: Visa, MasterCard, American Express, and Discover
`^((4\d{3})|(5[1-5]\d{2})|(6011))-?\d{4}-?\d{4}-?\d{4}|3[4,7]\d{13}$` (Continued)

| Regex | Description |
|--------|----------------------------------|
| [4,7] | A 4 or a 7, followed by |
| \d{13} | Any thirteen digits, followed by |
| \$ | End of line |

* **In English:** A number starting with 4 and three digits, or 5 and three digits, or 6011, followed by a hyphen, followed by three sets of four digits, or 34 and thirteen digits, or 37 and thirteen digits.

NOTE This regex does not, and cannot, conform to mod 10 verification. To find a Java program that does, please visit <http://www.influxs.com>.

Table C-7. Real Number `^[+-]?[d+](\.[d+])?$`

| Regex | Description |
|-------|---|
| ^ | Beginning of line, followed by |
| [+-]? | An optional plus or a minus sign |
| \d+ | Followed by one or more digits, followed by |
| (| A group consisting of |
| \. | A period, followed by |
| \d+ | One or more digits |
|)? | Close group, and make it optional |
| \$ | End of line |

* **In English:** Any number of digits followed by an optional decimal component.

Index

A

- addresses confirmation, 33–37
- alternates, 12–15
- append position, 99
- appendReplace method, 239
- appendReplacement method, 98–102, 123
- appendTail method, 98, 103, 239
- Apress Web site, 6
- ArrayList, 189

B

- back references, 20–21, 123–24
- batch reads and writes, 145–46
- boolean String.matches method, 42
- boundary characters, 11–12
- ByteBuffer, 144, 145
- ByteBuffers, 152–60

C

- Calendar object, 185
- CANON_EQ flag, 59, 230
- CASE_INSENSITIVE field, 57, 228
- character classes, 15–17
- characters
 - boundary, 11–12
 - common, 9–11
- CharSequence parameter, 56, 63
- Command object, 143

- COMMENT flag, 52
- COMMENTS flag, 57–58, 63, 228
- common characters, 9–11
- compile flags, 189
- compile method, 59–61, 230–32
- compiling patterns, 160–62
- composition technique, 7
- connections, optimizing, 144–45

D

- data validation, 42–46
- date confirmation, 27–30, 184–89
- delimiting strings, 106–7
- DOTALL flag, 58–59, 229
- duplicate words, finding, 38–41

E

- EDI document, validating, 207–9
- end method, 81–86, 237
- examineLog method, 168
- examples, 173–210
 - confirming date formats, 184–89
 - confirming format of phone number, 173–78
 - confirming zip codes, 179–83
 - extracting phone numbers from files, 200–203
 - modifying contents of files, 196–99

- searching directory for file containing expression, 203–7
- searching files, 192–96
- searching strings, 189–92
- validating EDI document, 207–9

F

- FileChannels, 144–45
 - modifying contents of files, 198
 - searching files, 192–94
 - storing patterns, 152–60
- FileInputStream, 144
- FileLocks tool, 144
- FileOutputStream, 144
- files
 - modifying contents of, 196–99
 - searching, 192–96
- find method, 93–96, 238
- finding duplicate words, 38–41
- find(int start) method, 238
- flags method, 56, 63–64, 232–33

G

- getFileContent method, 199
- getRegex method, 154
- getXML method, 151
- group method, 86–87, 237
 - and end method, 83, 84
 - and qualifiers, 127, 128
 - and start method, 79
 - and subgroups, 119, 120, 121
- groupCount method, 90–91, 238
- group(int group) method, 88–90, 237–38

- group(int) method, 88–89
- groups, 18–19, 117–18
 - and Matcher object, 71–72
 - noncapturing, 138, 142
 - subgroups, 119–22

I

- IllegalStateException, 93
 - and appendReplacement method, 102
 - and end method, 83, 86
 - and group method, 87, 90
 - and start method, 78, 81
- indexOf method, 113, 160
- IndexOutOfBoundsException, 81, 86, 90, 95, 102, 124

J

- Java
 - integrating with regular expressions, 21–41
 - confirming addresses example, 33–37
 - confirming dates example, 27–30
 - confirming phone number formats example, 22–25
 - confirming zip codes example, 25–27
 - finding duplicate words example, 38–41
 - overview, 21–22
 - vs. Perl, 51–52
- java.util.Properties, 173
- Java.util.regex object model, 55–115
 - Matcher object
 - appendReplacement method, 98–102
 - appendTail method, 103

- and back references, 123
- date validation with, 44–46
- end method, 81–86, 83–86
- find method, 93–96
- and finding duplicate words, 38–40
- group method, 86–87
- groupCount method, 90–91
- and groups, 71–72
- lookingAt method, 96–98
- matcher method, 63
- matches method, 91–92
- overview, 70–71
- pattern method, 73
- and Pattern.matches method, 64–65
- replaceAll method, 103–5
- replaceFirst method, 105–6
- reset method, 73–76
- start method, 77–81
- vs. String.matches method, 53
- and String.replaceAll method, 109
- and String.replaceFirst method, 108
- overview, 55
- Pattern object, 55–70
 - CANON_EQ flag, 59
 - CASE_INSENSITIVE field, 57
 - COMMENTS flag, 57–58
 - compile method, 59–61
 - DOTALL flag, 58–59
 - flags method, 63–64
 - matcher method, 63
 - matches method, 64–65
 - MULTILINE flag, 58

- overview, 55–56
- pattern method, 62–63
- split method, 65–70
- UNICODE_CASE flag, 59
- UNIX_LINES flag, 57
- String object, 106–12
 - delimiting strings, 106–7
 - matches method, 107–8
 - replaceFirst method, 108–9
 - split method, 109–12

L

- limit parameter, 68, 110
- LinkedHashMap, 191
- Logger.throwing method, 162
- lookarounds, 130–37
 - negative lookaheads, 132–34
 - negative lookbehinds, 137
 - overview, 130
 - positive lookaheads, 130–32
 - positive lookbehinds, 134–37
- lookingAt method, 96–98

M

- Map, 191
- MappedByteBuffer, 195
- matcher method, 63, 232
- Matcher object
 - appendReplacement method, 98–102
 - appendTail method, 103
 - and back references, 123
 - date validation with, 44–46

- end method, 81–86, 83–86
- find method, 93–96
 - and finding duplicate words, 38–40
- group method, 86–87
- groupCount method, 90–91
- and groups, 71–72
- lookingAt method, 96–98
- matcher method, 63
- matches method, 91–92
- overview, 70–71
- pattern method, 73
 - and Pattern.matches method, 64–65
- replaceAll method, 103–5
- replaceFirst method, 105–6
- reset method, 73–76
- start method, 77–81
- vs. String.matches method, 53
 - and String.replaceAll method, 109
 - and String.replaceFirst method, 108
- matcher(CharSequence input) method, 63
- matchers, defining, 3–4
- matches
 - specifying position of, 139
 - specifying size of, 140
- matches method, 53, 91–92, 238
- matches(String regex, CharSequence input) method, 64
- matches(String regex,CharSequence input) method, 64–65, 233–34
- matches(String regex) method, 107–8
- MatchNameFormats.java program, 30
- metacharacters, 5
- Method class methods, 236–39

modifying contents of files, 196–99

MULTILINE flag, 58, 228–29

N

negation example, 17

negative lookaheads, 132–34

negative lookbehinds, 137

newsgroups, 210

noncapturing groups, 138, 142

NullPointerException, 108, 109, 110, 112

O

object-oriented regular expressions,
143–71

- batch reads and writes, 145–46

- compiling patterns as needed, 160–62

- and not limiting to regex solutions,
163–68

- optimizing connections, 144–45

- overview, 143

- storing patterns externally, 146–60

- not using normal property file,
147–48

- not using XML, 148–52

- overview, 146–47

- using FileChannels and ByteBuffers,
152–60

optimizing connections, 144–45

optimizing regular expressions, 138–40

P

Pattern class fields, 227–30

Pattern class methods, 230–36

pattern method, 62–63, 73, 232, 236

- Pattern object, 55–70
 - and appendReplacement method, 98
 - CANON_EQ flag, 59
 - CASE_INSENSITIVE field, 57
 - COMMENTS flag, 57–58
 - compile method, 59–61
 - data validation with, 44–46
 - DOTALL flag, 58–59
 - finding duplicate words example, 38–40
 - flags method, 63–64
 - and lookingAt method, 96
 - matcher method, 63, 232
 - matches method, 64–65, 108
 - MULTILINE flag, 58
 - overview, 55–56
 - pattern method, 62–63
 - and reluctant qualifiers, 130
 - split method, 65–70
 - split(String regex, int limit) method, 112
 - split(String regex) method, 110
 - vs. String.matches method, 53
 - and String.replaceFirst method, 109
 - UNICODE_CASE flag, 59
 - UNIX_LINES flag, 57
 - patterns
 - compiling as needed, 160–62
 - creating, 5–7
 - defining, 2–3
 - reading, 8–9
 - storing externally, 146–60
 - not using normal property file, 147–48
 - not using XML, 148–52
 - overview, 146–47
 - using FileChannels and ByteBuffers, 152–60
 - PatternSyntaxException, 55, 65, 108, 109, 110, 112, 230–31
 - Perl, vs. Java's regex support, 51–52
 - phone numbers
 - extracting from files, 200–203
 - format confirmation, 22–25, 173–78
 - positive lookaheads, 130–32
 - positive lookbehinds, 134–37
 - Properties object, 147, 148, 152
 - pull technique, 5–6
 - push technique, 6–7
- ## Q
- qualifiers, 125–30
 - quantifiers and alternates, 12–15
- ## R
- reading patterns, 8–9
 - RegexProperties class, 173, 179
 - regular expression operations, 41–51
 - data validation, 42–46
 - overview, 41
 - search and replace, 46–48
 - splitting strings, 48–51
 - regular expressions, building blocks of, 2–4
 - replaceAll method, 46–48, 109, 115, 123, 239
 - replaceAll(String replacement) method, 103–5
 - replaceFirst method, 46–48, 123, 239

- replaceFirst(String regex,String replacement) method, 108–9
- replaceFirst(String replacement) method, 105–6
- reset method, 73–76, 106, 236
- resources, 210
- ResultSet, 3, 117, 143
- RuntimeException, 60, 62
- RX.java program, 6
- S**
- saveXML method, 151
- searchFile method, 194
- searching, 46–48
 - directory for file containing expression, 203–7
 - files, 192–96
 - strings, 189–92
- searchString method, 192
- split(CharSequence input, int limit) method, 68–70, 235–36
- split(CharSequence input) method, 65–68, 234–35
- split(String regex, int limit) method, 110–12
- split(String regex) method, 109–10
- splitting strings, 48–51
- start method, 77–78, 97
 - and find method, 93
 - and group method, 87
 - and replaceFirst method, 106
- start(int group) method, 78–81, 237

- storing patterns externally, 146–60
 - not using normal property file, 147–48
 - not using XML, 148–52
 - overview, 146–47
 - using FileChannels and ByteBuffers, 152–60
- String object, 106–12
 - and common characters, 11
 - delimiting strings, 106–7
 - group method, 86–90
 - and groups, 71, 72
 - indexOf method, 113, 160
 - matcher method, 43
 - matches method, 15, 53
 - matches(regex) method, 38
 - matches(String regex,CharSequence input) method, 64
 - matches(String regex) method, 22, 45, 107–8
 - modification of, 52
 - replaceAll method, 103–5, 109, 175
 - replaceFirst method, 108, 108–9
 - split method, 68, 113, 114, 185
 - split(CharSequence input, int limit) method, 68–70
 - split(CharSequence input) method, 65–68
 - split(String regex, int limit) method, 110–12
 - split(String regex) method, 109–10
 - substring method, 46, 87

StringBuffer object, 98, 99, 100–101, 103

strings

- data validation with, 42–44

- prechecking candidate strings, 138

- searching, 189–92

- splitting, 48–51

StringTokenizer, 1, 109, 163

subgroups, 119–22

substrings method, 46

syntax of regular expressions, 7–20

- back references, 20–21

- boundary characters, 11–12

- character classes, 15–17

- common characters, 9–11

- groups, 18–19

- overview, 7

- quantifiers and alternates, 12–15

- reading patterns, 8–9

T

telephone numbers. *See* phone numbers

U

UNICODE_CASE flag, 59, 230

UNIX_LINES flag, 57, 227

V

validating EDI document, 207–9

W

words, duplicate, 38–41

writing regular expressions, 5–7

X

XML, not using for storing patterns, 148–52

XMLDecoder class, 148

XMLDecoder object, 149–50

XMLEncoder class, 148

XMLEncoder object, 149–50

XMLHelper class, 151

Z

zip code confirmation, 25–27, 179–83

zipPatternKey, 182

ASPToday

ASPToday is a unique solutions library for professional ASP Developers, giving quick and convenient access to a constantly growing library of **over 1000 practical and relevant articles and case studies**. We aim to publish a completely original professionally written and reviewed article every working day of the year. Consequently our resource is completely without parallel in the industry. Thousands of web developers use and recommend this site for real solutions, keeping up to date with new technologies, or simply increasing their knowledge.

Exciting Site Features!

Find it FAST!

Powerful full-text search engine so you can find exactly the solution you need.

Printer-friendly!

Print articles for a bound archive and quick desk reference.

Working Sample Code Solutions!

Many articles include complete downloadable sample code ready to adapt for your own projects.

ASPToday covers a broad range of topics including:

- | | |
|-----------------------|-----------------------|
| ▶ ASP.NET 1.x and 2.0 | ▶ Security |
| ▶ ADO.NET and SQL | ▶ Site Design |
| ▶ XML | ▶ Site Admin |
| ▶ Web Services | ▶ SMTP and Mail |
| ▶ E-Commerce | ▶ Classic ASP and ADO |

and much, much more...

To receive a **FREE** two-month subscription to ASPToday, visit
www.asptoday.com/subscribe.aspx and answer the question about this book!

The above **FREE** two-month subscription offer is good for six months from original copyright date of book this ad appears in. Each book will require a different promotional code to get this free offer- this code will determine the offer expiry date. Paid subscribers to ASPToday will receive 50% off of selected Apress books with a paid 3-month or one-year subscription. Subscribers will also receive discount offers and promotional email from Apress unless their subscriber preferences indicate they don't wish this. Offer limited to one **FREE** two-month subscription offer per person.



forums.apress.com

FOR PROFESSIONALS BY PROFESSIONALS ^{11.1}

JOIN THE APRESS FORUMS AND BE PART OF OUR COMMUNITY. You'll find discussions that cover topics of interest to IT professionals, programmers, and enthusiasts just like you. If you post a query to one of our forums, you can expect that some of the best minds in the business—especially Apress authors, who all write with *The Expert's Voice*™—will chime in to help you. Why not aim to become one of our most valuable participants (MVPs) and win cool stuff? Here's a sampling of what you'll find:

DATABASES

Data drives everything.

Share information, exchange ideas, and discuss any database programming or administration issues.

PROGRAMMING/BUSINESS

Unfortunately, it is.

Talk about the Apress line of books that cover software methodology, best practices, and how programmers interact with the "suits."

INTERNET TECHNOLOGIES AND NETWORKING

Try living without plumbing (and eventually IPv6).

Talk about networking topics including protocols, design, administration, wireless, wired, storage, backup, certifications, trends, and new technologies.

WEB DEVELOPMENT/DESIGN

Ugly doesn't cut it anymore, and CGI is absurd.

Help is in sight for your site. Find design solutions for your projects and get ideas for building an interactive Web site.

JAVA

We've come a long way from the old Oak tree.

Hang out and discuss Java in whatever flavor you choose: J2SE, J2EE, J2ME, Jakarta, and so on.

SECURITY

Lots of bad guys out there—the good guys need help.

Discuss computer and network security issues here. Just don't let anyone else know the answers!

MAC OS X

All about the Zen of OS X.

OS X is both the present and the future for Mac apps. Make suggestions, offer up ideas, or boast about your new hardware.

TECHNOLOGY IN ACTION

Cool things. Fun things.

It's after hours. It's time to play. Whether you're into LEGO® MINDSTORMS™ or turning an old PC into a DVR, this is where technology turns into fun.

OPEN SOURCE

Source code is good; understanding (open) source is better.

Discuss open source technologies and related topics such as PHP, MySQL, Linux, Perl, Apache, Python, and more.

WINDOWS

No defenestration here.

Ask questions about all aspects of Windows programming, get help on Microsoft technologies covered in Apress books, or provide feedback on any Apress Windows book.

HOW TO PARTICIPATE:

Go to the Apress Forums site at <http://forums.apress.com/>.

Click the New User link.