

# regexfilter.java File Reference

This application implements a text filter for multiple replacements using regular expressions (regex).

General purpose filter, replaces regular expressions (regex) in the input stream, in block or line modes. It works like egrep and sed, or like Perl regex, but it is simpler and it can be used on Unix/Linux and Windows O.S.

## Precondition:

- input: standard input (allows piping) or file.
- output: standard output or file
- regex: all pairs regex/replacement are in a file.
- works line by line or on whole input.
- macro expansions in replacement strings.

## Author:

M. Sillano ([marco.sillano@gmail.com](mailto:marco.sillano@gmail.com))

## Version:

4.02 10/11/25 (c) M.Sillano 2006-2011

Filters the input stream, via `java.String.replaceAll()`, using regex/replacement pairs.

All regex/replacement pairs are stored in an ASCII file, and they are all processed in sequence for every input line or for the whole input in block mode.

## use

```
Usage:      regexfilter -h|-?|--help|--version",
            regexfilter [-bx] [datFile] [-i=FILE] [-u=FILE]",
```

Filters the inputFile, using regex/replacement pairs in datFile.

```
options:  -h|-?|--help    display this help and exit.
           --version      print version and exit.
           -b             block mode. Default = line mode.
           -x             datFile in XML mode. Default = plain text.
-i=FILE,  --input=FILE    the text input file.
                        Default = standard input.
-o=FILE,  --output=FILE   the text output file.
                        Default = standard output.
datFile:  regex/replacement pairs file. Default = ./regexfilter.dat.
           if not found, builds an example file.
```

## datFile - syntax:

- regex must be double-escaped ("\\") using regex rules (not in XML mode)
- in replacement some chars must be single-escaped ( \", \', \\, \t, \f, \n, \r [for: ", ', \, TAB, FF, LF, CR] ) and same for hexadecimal unicode: \ uxxxx. Some must be double-escaped (\\space, \\\$ [for : initial-final spaces, \$ char]).
- in replacement \$n is a references to 'n' captured subsequences.
- replacement allows empty lines (delete)

- macros in replacement (processed before regexp):
  - `\\space` is for initial/final spaces
  - `$V` is regexfilter version
  - `$F` is full path of the input file
  - `$P` is the path (file excluded) of the input file
  - `$N` is the file name (no ext.) of the input file
  - `$E` is the extension (no dot) of the input file
  - `$D` is the filter (file DAT) name (no ext.)
  - `$C` is the Copyright notice (CC BY-SA)
  - `$T` is the timestamp

## datFile - text format

```
# rule 1: replaces "# /**" with "/*"
regex1=^#[ ]*(/\\*\\*)
replacement1=$1

# rule 302: replaces "# *" with "*"
regex302=^#[ ]*\\*
replacement302=*
```

*note: regex/replacement indexes MUST be in 1..999 interval.*

## datFile - XML format

```
<comment> rule 1: replaces "# /**" with "/*" </comment>
<entry key="regex1">^#[ ]*(/\\*\\*)</entry>
<entry key="replacement1">$1</entry>
<comment> rule 2: replaces "# *" with "*" </comment>
<entry key="regex2">^#[ ]*\\*</entry>
<entry key="replacement2">*</entry>
```

*note: regex/replacement indexes MUST start from 1 and MUST be all consecutive.*

## Summary of regular-expression constructs

(from JDK 6 documentation)

Construct	Matches
<b>Characters</b>	
<code>x</code>	The character <code>x</code>
<code>\\</code>	The backslash character
<code>\\0n</code>	The character with octal value <code>0n</code> ( $0 \leq n \leq 7$ )
<code>\\0nn</code>	The character with octal value <code>0nn</code> ( $0 \leq n \leq 7$ )
<code>\\0mnn</code>	The character with octal value <code>0mnn</code> ( $0 \leq m \leq 3, 0 \leq n \leq 7$ )
<code>\\xhh</code>	The character with hexadecimal value <code>0xhh</code>
<code>\\uhhhh</code>	The character with hexadecimal value <code>0xhhhh</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> )

`\cx` The control character corresponding to *x*

### Character classes

`[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

`.` 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\r\b\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)

`\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}\x20]`  
`\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\r\b\f\r]`

### java.lang.Character classes (simple java character type)

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

### Classes for Unicode blocks and categories

`\p{InGreek}` A character in the Greek block (simple bloc)  
`\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

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

### Greedy quantifiers

<code>X?</code>	$X$ , once or not at all
<code>X*</code>	$X$ , zero or more times
<code>X+</code>	$X$ , one or more times
<code>X{n}</code>	$X$ , exactly $n$ times
<code>X{n,}</code>	$X$ , at least $n$ times
<code>X{n,m}</code>	$X$ , at least $n$ but not more than $m$ times

### Reluctant quantifiers

<code>X??</code>	$X$ , once or not at all
<code>X*?</code>	$X$ , zero or more times
<code>X+?</code>	$X$ , one or more times
<code>X{n}?</code>	$X$ , exactly $n$ times
<code>X{n,}?</code>	$X$ , at least $n$ times
<code>X{n,m}?</code>	$X$ , at least $n$ but not more than $m$ times

### Possessive quantifiers

<code>X?+</code>	$X$ , once or not at all
<code>X*+</code>	$X$ , zero or more times
<code>X++</code>	$X$ , one or more times
<code>X{n}+</code>	$X$ , exactly $n$ times
<code>X{n,}+</code>	$X$ , at least $n$ times
<code>X{n,m}+</code>	$X$ , at least $n$ but not more than $m$ times

### Logical operators

<code>XY</code>	$X$ followed by $Y$
<code>X Y</code>	Either $X$ or $Y$
<code>(X)</code>	$X$ , as a capturing group

### Back references

<code>\n</code>	Whatever the $n^{\text{th}}$ capturing group matched
-----------------	--

### Quotation

<code>\</code>	Nothing, but quotes the metacharacters: <code>&lt;([{\^-= \$!  ]})? * + . &gt;</code>
<code>\Q</code>	Nothing, but quotes all characters until <code>\E</code>
<code>\E</code>	Nothing, but ends quoting started by <code>\Q</code>

### Special constructs (non-capturing)

<code>(?:X)</code>	$X$ , as a non-capturing group
<code>(?idsux-idmsux)</code>	Nothing, but turns match flags <code>i d m s u x</code> on - off
<code>(?idsux-idmsux:X)</code>	$X$ , as a non-capturing group with the given flags <code>i d m s u x</code> on - off
<code>(?=X)</code>	$X$ , via zero-width positive lookahead
<code>(?!X)</code>	$X$ , via zero-width negative lookahead
<code>(?&lt;=X)</code>	$X$ , via zero-width positive lookbehind
<code>(?&lt;!X)</code>	$X$ , via zero-width negative lookbehind
<code>(?&gt;X)</code>	$X$ , as an independent, non-capturing group

Copyright 2008 Sun Microsystems, Inc. Reprinted with permission (see <http://java.sun.com/javase/6/docs/api/java/util/regex/Pattern.html#sum>)

## line terminators

The following are recognized as line terminators in regex:

- A newline (line feed) character ("`\n`"), - UNIX
- A carriage-return character followed immediately by a newline character ("`\r\n`"), - DOS/WIN
- A standalone carriage-return character ("`\r`"), - MAC
- A next-line character ("`\u0085`"),
- A line-separator character ("`\u2028`")
- A paragraph-separator character ("`\u2029`").

## flag expressions (?idmsux-idmsux)

- In *block mode* [-b] the expressions "`^`" and "`$`" match at the beginning and the end of the entire input sequence. The embedded flag expression "(?m)" enables multiline mode, so these expressions match just after or just before, respectively, any line terminator.
- In *dotall mode*, the expression "`.`" matches any character, including a line terminator. By default this expression does not match line terminators. Dotall mode can be enabled via the embedded flag expression "(?s)".
- *Case-insensitive* matching for ASCII can be enabled via the embedded "(?i)".
- *Unicode-aware* case folding can be enabled via the embedded flag expression "(?u)". Requires "(?i)".
- *Comments mode* can be enabled via the embedded flag expression "(?x)". In this mode, whitespace is ignored, and embedded comments starting with `#` are ignored until the end of a line.
- *Unix-lines* mode can be enabled via the embedded flag expression "(?d)". In this mode only "`\n`" line terminator is recognized in the behavior of "`.`", "`^`", and "`$`".