

Need to Know

Pattern arguments in stringr are interpreted as regular expressions *after any special characters have been parsed*.

In R, you write regular expressions as *strings*, sequences of characters surrounded by quotes (""), or single quotes ('').

Some characters cannot be represented directly in an R string. These must be represented as **special characters**, sequences of characters that have a specific meaning, e.g.

Special Character Represents

```
\\      \"      new line
\\n     \"      new line
Run ?"" to see a complete list
```

Because of this, whenever a \ appears in a regular expression, you must write it as \\ in the string that represents the regular expression.

Use **writeLines()** to see how R views your string after all special characters have been parsed.

```
writeLines("\\\\")
# \\
writeLines("\\\\ is a backslash")
# \\ is a backslash
```

INTERPRETATION

Patterns in stringr are interpreted as regexes. To change this default, wrap the pattern in one of:

regex() (pattern, ignore.case = FALSE, multiline = FALSE, comments = FALSE, dotall = FALSE, ...) Modifies a regex to ignore cases, match end of lines as well as end of strings, allow R comments within regex's, and/or to have . match everything including \n.
str_detect("hi", regex("i", TRUE))

fixed() Matches raw bytes but will miss some characters that can be represented in multiple ways (fast). str_detect("u0130", fixed("ı"))

coll() Matches raw bytes and will use locale specific collation rules to recognize characters that can be represented in multiple ways (slow). str_detect("u0130", coll("i", TRUE, locale = "tr"))

boundary() Matches boundaries between characters, line_breaks, sentences, or words. str_split(sentences, boundary("word"))

Regular Expressions -

Regular expressions, or *regexps*, are a concise language for describing patterns in strings.

MATCH CHARACTERS

string (type this)	regex (to mean this)	matches (which matches this)	example
<code>a</code>	<code>(etc.)</code>	<code>a</code> (etc.)	see("a")
<code>.</code>	<code>(etc.)</code>	<code>.</code>	see("\\.")
<code>!</code>	<code>(etc.)</code>	<code>!</code>	see("\\!")
<code>?</code>	<code>(etc.)</code>	<code>?</code>	see("\\?")
<code>\\</code>	<code>(etc.)</code>	<code>\\</code>	see("\\\\")
<code>(</code>	<code>(etc.)</code>	<code>(</code>	see("\\(")
<code>)</code>	<code>(etc.)</code>	<code>)</code>	see("\\)")
<code>{</code>	<code>(etc.)</code>	<code>{</code>	see("\\{")
<code>}</code>	<code>(etc.)</code>	<code>}</code>	see("\\}")
<code>\\n</code>	<code>(etc.)</code>	new line (return)	see("\\n")
<code>\\t</code>	<code>(etc.)</code>	tab	see("\\t")
<code>\\s</code>	<code>(etc.)</code>	any whitespace (\\S for non-whitespaces)	see("\\s")
<code>\\d</code>	<code>(etc.)</code>	any digit (\\D for non-digits)	see("\\d")
<code>\\w</code>	<code>(etc.)</code>	any word character (\\W for non-word chars)	see("\\w")
<code>\\b</code>	<code>(etc.)</code>	word boundaries	see("\\b")
<code>[digit:]</code>	<code>(etc.)</code>	digits	see("[digit:]")
<code>[alpha:]</code>	<code>(etc.)</code>	letters	see("[alpha:]")
<code>[lower:]</code>	<code>(etc.)</code>	lowercase letters	see("[lower:]")
<code>[upper:]</code>	<code>(etc.)</code>	uppercase letters	see("[upper:]")
<code>[alnum:]</code>	<code>(etc.)</code>	letters and numbers	see("[alnum:]")
<code>[punct:]</code>	<code>(etc.)</code>	punctuation	see("[punct:]")
<code>[graph:]</code>	<code>(etc.)</code>	letters, numbers, and punctuation	see("[graph:]")
<code>[space:]</code>	<code>(etc.)</code>	space characters (i.e. \\s)	see("[space:]")
<code>[blank:]</code>	<code>(etc.)</code>	space and tab (but not new line)	see("[blank:]")
<code>.</code>	<code>(etc.)</code>	every character except a new line	see(".")

¹ Many base R functions require classes to be wrapped in a second set of [], e.g. `[[:digit:]]`

ALTERNATES

regex	matches	example
<code>ab d</code>	or	alt("ab d")
<code>[abe]</code>	one of	alt("[abe]")
<code>[^abe]</code>	anything but	alt("[^abe]")
<code>[a-c]</code>	range	alt("[a-c]")

ANCHORS

regex	matches	example
<code>^</code>	start of string	anchor("^a")
<code>\$</code>	end of string	anchor("a\$")

LOOK AROUNDS

regex	matches	example
<code>a(?=)</code>	followed by	look("a(?=c)")
<code>a(?!)</code>	not followed by	look("a(?!c)")
<code>(?<=)</code>	preceded by	look("(?<=b)a")
<code>(?<!)</code>	not preceded by	look("(?<!b)a")

see <- function(rx) str_view_all("abc ABC 123 t: !? \\ \\n", rx)

string	regex	matches	example
abc ABC 123	<code>!?</code>	! ?	see("a")
abc ABC 123	<code>!?</code>	! ?	see("\\.")
abc ABC 123	<code>!?</code>	! ?	see("\\!")
abc ABC 123	<code>!?</code>	! ?	see("\\?")
abc ABC 123	<code>!?</code>	! ?	see("\\\\")
abc ABC 123	<code>!?</code>	! ?	see("\\(")
abc ABC 123	<code>!?</code>	! ?	see("\\)")
abc ABC 123	<code>!?</code>	! ?	see("\\{")
abc ABC 123	<code>!?</code>	! ?	see("\\}")
abc ABC 123	<code>!?</code>	! ?	see("\\n")
abc ABC 123	<code>!?</code>	! ?	see("\\t")
abc ABC 123	<code>!?</code>	! ?	see("\\s")
abc ABC 123	<code>!?</code>	! ?	see("\\d")
abc ABC 123	<code>!?</code>	! ?	see("\\w")
abc ABC 123	<code>!?</code>	! ?	see("\\b")
abc ABC 123	<code>!?</code>	! ?	see("[digit:]")
abc ABC 123	<code>!?</code>	! ?	see("[alpha:]")
abc ABC 123	<code>!?</code>	! ?	see("[lower:]")
abc ABC 123	<code>!?</code>	! ?	see("[upper:]")
abc ABC 123	<code>!?</code>	! ?	see("[alnum:]")
abc ABC 123	<code>!?</code>	! ?	see("[punct:]")
abc ABC 123	<code>!?</code>	! ?	see("[graph:]")
abc ABC 123	<code>!?</code>	! ?	see("[space:]")
abc ABC 123	<code>!?</code>	! ?	see("[blank:]")
abc ABC 123	<code>!?</code>	! ?	see(".")

QUANTIFIERS

regex	matches	example
<code>a?</code>	zero or one	quant("a?")
<code>a*</code>	zero or more	quant("a*")
<code>a+</code>	one or more	quant("a+")
<code>a{n}</code>	exactly n	quant("a{2}")
<code>a{n,}</code>	n or more	quant("a{2,}")
<code>a{n,m}</code>	between n and m	quant("a{2,4}")

GROUPS

regex	matches	example
<code>(ab d)e</code>	sets precedence	alt("(ab d)e")
<code>string</code>	(type this)	(the result is the same as ref("abba"))
<code>\\1</code>	first () group, etc.	ref("(a)(b)(2)(1)")



[:space:]	new line
↵	
[:blank:]	space
␣	
[:graph:]	
· , ; : ! \ / ' = * + - ^	
[:punct:]	
~ " ' [] { } () < > @ # \$	
[:alnum:]	
0 1 2 3 4 5 6 7 8 9	
[:digit:]	
[:alpha:]	
a b c d e f A B C D E F	
g h i j k l G H I J K L	
m n o p q r M N O P Q R	
s t u v w x S T U V W X	
z Z	

quant <- function(rx) str_view_all("a.aa.aaa", rx)

ref <- function(rx) str_view_all("abbaab", rx)