

# 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
\\	\
\"	"
\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 regexs. 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 of end of strings, allow R comments within regex's, and/or to have . match everything including \n.  
`str_detect("i", regex("i", TRUE))`

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

**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"))`

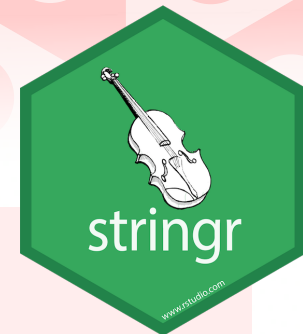
# Regular Expressions - Regular expressions, or *regexps*, are a concise language for describing patterns in strings.

## MATCH CHARACTERS

string (type this)	regex (to mean this)	matches (which matches this)	example
	<b>a (etc.)</b>	a (etc.)	<code>see("a")</code> abc ABC 123 .!?\()\
\\.	\\.	.	<code>see("\\.")</code> abc ABC 123 .!?\()\
\\!	\\!	!	<code>see("\\!")</code> abc ABC 123 .!?\()\
\\?	\\?	?	<code>see("\\?")</code> abc ABC 123 .!?\()\
\\\\	\\\\	\\	<code>see("\\\\")</code> abc ABC 123 .!?\()\
\\(	\\(	(	<code>see("\\(")</code> abc ABC 123 .!?\()\
\\)	\\)	)	<code>see("\\)")</code> abc ABC 123 .!?\()\
\\{	\\{	{	<code>see("\\{")</code> abc ABC 123 .!?\()\
\\}	\\}	}	<code>see("\\}")</code> abc ABC 123 .!?\()\
\\n	\\n	new line (return)	<code>see("\\n")</code> abc ABC 123 .!?\()\
\\t	\\t	tab	<code>see("\\t")</code> abc ABC 123 .!?\()\
\\s	\\s	any whitespace ( <b>S</b> for <i>non-whitespaces</i> )	<code>see("\\s")</code> abc ABC 123 .!?\()\
\\d	\\d	any digit ( <b>D</b> for <i>non-digits</i> )	<code>see("\\d")</code> abc ABC 123 .!?\()\
\\w	\\w	any word character ( <b>W</b> for <i>non-word chars</i> )	<code>see("\\w")</code> abc ABC 123 .!?\()\
\\b	\\b	word boundaries	<code>see("\\b")</code> abc ABC 123 .!?\()\
	<b>[digit:]</b> <sup>1</sup>	digits	<code>see("[digit:]")</code> abc ABC 123 .!?\()\
	<b>[alpha:]</b> <sup>1</sup>	letters	<code>see("[alpha:]")</code> abc ABC 123 .!?\()\
	<b>[lower:]</b> <sup>1</sup>	lowercase letters	<code>see("[lower:]")</code> abc ABC 123 .!?\()\
	<b>[upper:]</b> <sup>1</sup>	uppercase letters	<code>see("[upper:]")</code> abc ABC 123 .!?\()\
	<b>[alnum:]</b> <sup>1</sup>	letters and numbers	<code>see("[alnum:]")</code> abc ABC 123 .!?\()\
	<b>[punct:]</b> <sup>1</sup>	punctuation	<code>see("[punct:]")</code> abc ABC 123 .!?\()\
	<b>[graph:]</b> <sup>1</sup>	letters, numbers, and punctuation	<code>see("[graph:]")</code> abc ABC 123 .!?\()\
	<b>[space:]</b> <sup>1</sup>	space characters (i.e. \s)	<code>see("[space:]")</code> abc ABC 123 .!?\()\
	<b>[blank:]</b> <sup>1</sup>	space and tab (but not new line)	<code>see("[blank:]")</code> abc ABC 123 .!?\()\
	.	every character except a new line	<code>see(".")</code> abc ABC 123 .!?\()\

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

<b>[space:]</b>
↵ new line
<b>[blank:]</b>
□ space
□ tab



graph

punct

. , : ; ? ! / \* @ #  
- \_ " ' [ ] { } ( )

symbol

| ` = + ^  
~ < > \$

alnum

digit

0 1 2 3 4 5 6 7 8 9

alpha

lower

a b c d e f  
g h i j k l  
m n o p q r  
s t u v w x  
y z

upper

A B C D E F  
G H I J K L  
M N O P Q R  
S T U V W X  
Y Z

## ALTERNATES

`alt <- function(rx) str_view("abcde", rx)`

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

## ANCHORS

`anchor <- function(rx) str_view("aaa", rx)`

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

## LOOK AROUNDS

`look <- function(rx) str_view("bacad", rx)`

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

## QUANTIFIERS

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

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

## GROUPS

`ref <- function(rx) str_view("abbaab", rx)`

Use parentheses to set precedent (order of evaluation) and create groups

regex	matches	example
<b>(ab d)e</b>	sets precedence	<code>alt("(ab d)e")</code> abcde

Use an escaped number to refer to and duplicate parentheses groups that occur earlier in a pattern. Refer to each group by its order of appearance

string (type this)	regex (to mean this)	matches (which matches this)	example (the result is the same as ref("abba"))
<b>\\1</b>	<b>\\1</b> (etc.)	first () group, etc.	<code>ref("(a)(b)\\2\\1")</code> abbaab