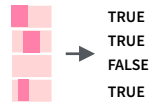


# String manipulation with stringr : : CHEAT SHEET

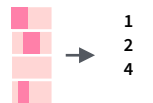


The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.

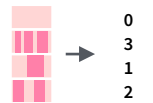
## Detect Matches



**str\_detect**(string, **pattern**) Detect the presence of a pattern match in a string. `str_detect(fruit, "a")`



**str\_which**(string, **pattern**) Find the indexes of strings that contain a pattern match. `str_which(fruit, "a")`



**str\_count**(string, **pattern**) Count the number of matches in a string. `str_count(fruit, "a")`

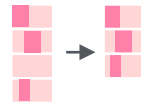


**str\_locate**(string, **pattern**) Locate the positions of pattern matches in a string. Also **str\_locate\_all**. `str_locate(fruit, "a")`

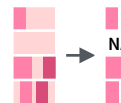
## Subset Strings



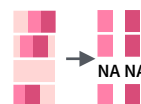
**str\_sub**(string, start = 1L, end = -1L) Extract substrings from a character vector. `str_sub(fruit, 1, 3); str_sub(fruit, -2)`



**str\_subset**(string, **pattern**) Return only the strings that contain a pattern match. `str_subset(fruit, "b")`

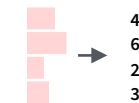


**str\_extract**(string, **pattern**) Return the first pattern match found in each string, as a vector. Also **str\_extract\_all** to return every pattern match. `str_extract(fruit, "[aeiou]")`

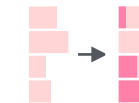


**str\_match**(string, **pattern**) Return the first pattern match found in each string, as a matrix with a column for each ( ) group in pattern. Also **str\_match\_all**. `str_match(sentences, "(a|the) ([^ ]+)")`

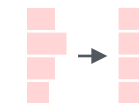
## Manage Lengths



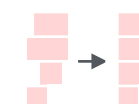
**str\_length**(string) The width of strings (i.e. number of code points, which generally equals the number of characters). `str_length(fruit)`



**str\_pad**(string, width, side = c("left", "right", "both"), pad = " ") Pad strings to constant width. `str_pad(fruit, 17)`



**str\_trunc**(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. `str_trunc(fruit, 3)`

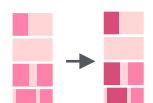


**str\_trim**(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. `str_trim(fruit)`

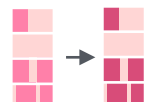
## Mutate Strings



**str\_sub()** <- value. Replace substrings by identifying the substrings with `str_sub()` and assigning into the results. `str_sub(fruit, 1, 3) <- "str"`



**str\_replace**(string, **pattern**, replacement) Replace the first matched pattern in each string. `str_replace(fruit, "a", "-")`



**str\_replace\_all**(string, **pattern**, replacement) Replace all matched patterns in each string. `str_replace_all(fruit, "a", "-")`

A STRING  
↓  
a string

**str\_to\_lower**(string, locale = "en")<sup>1</sup> Convert strings to lower case. `str_to_lower(sentences)`

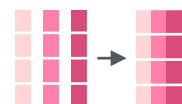
a string  
↓  
A STRING

**str\_to\_upper**(string, locale = "en")<sup>1</sup> Convert strings to upper case. `str_to_upper(sentences)`

a string  
↓  
A String

**str\_to\_title**(string, locale = "en")<sup>1</sup> Convert strings to title case. `str_to_title(sentences)`

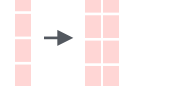
## Join and Split



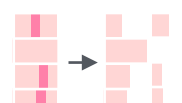
**str\_c**(..., sep = "", collapse = NULL) Join multiple strings into a single string. `str_c(letters, LETTERS)`



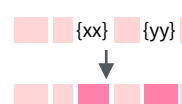
**str\_c**(..., sep = "", collapse = "") Collapse a vector of strings into a single string. `str_c(letters, collapse = "")`



**str\_dup**(string, times) Repeat strings times times. `str_dup(fruit, times = 2)`



**str\_split\_fixed**(string, **pattern**, n) Split a vector of strings into a matrix of substrings (splitting at occurrences of a pattern match). Also **str\_split** to return a list of substrings. `str_split_fixed(fruit, " ", n=2)`



**str\_glue**(..., .sep = "", .envir = parent.frame()) Create a string from strings and {expressions} to evaluate. `str_glue("Pi is {pi}")`



**str\_glue\_data**(.x, ..., .sep = "", .envir = parent.frame(), .na = "NA") Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. `str_glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")`

## Order Strings



**str\_order**(x, decreasing = FALSE, na\_last = TRUE, locale = "en", numeric = FALSE, ...) <sup>1</sup> Return the vector of indexes that sorts a character vector. `x[str_order(x)]`



**str\_sort**(x, decreasing = FALSE, na\_last = TRUE, locale = "en", numeric = FALSE, ...) <sup>1</sup> Sort a character vector. `str_sort(x)`

## Helpers

apple  
banana  
pear

**str\_conv**(string, encoding) Override the encoding of a string. `str_conv(fruit, "ISO-8859-1")`

apple  
banana  
pear

**str\_view**(string, **pattern**, match = NA) View HTML rendering of first regex match in each string. `str_view(fruit, "[aeiou]")`

**str\_view\_all**(string, **pattern**, match = NA) View HTML rendering of all regex matches. `str_view_all(fruit, "[aeiou]")`

**str\_wrap**(string, width = 80, indent = 0, exdent = 0) Wrap strings into nicely formatted paragraphs. `str_wrap(sentences, 20)`

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

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

**[space:]**  
new line

**[blank:]**  
space  
tab

**[graph:]**

**[punct:]**

. , : ; ? ! \ | / ` = \* + - ^  
\_ ~ " ' [ ] { } ( ) < > @ # \$

**[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  
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  
Z

## ALTERNATES

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

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

## ANCHORS

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

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

## LOOK AROUNDS

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

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

## QUANTIFIERS

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

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

## GROUPS

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

regex	matches	example
<b>(ab d)e</b>	sets precedence	alt("(ab d)e") 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"))
<code>\\1</code>	<code>\\1 (etc.)</code>	first () group, etc.	ref("(a)(b)\\2\\1") abbaab

