# 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

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

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

In the example column, each row shows: abc ABC 123 .!?\(){}

## [:space:]
↵ new line

## [:blank:]   .
□ space
□ tab

## [:graph:]

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

### [:alnum:]

#### [:digit:]
0 1 2 3 4 5 6 7 8 9

#### [:alpha:]

| [:lower:] | [:upper:] |
|---|---|
| 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 |

---

## ALTERNATES

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

| regexp | matches | example | |
|---|---|---|---|
| ab\|d | or | alt("ab\|d") | abcde |
| [abe] | one of | alt("[abe]") | abcde |
| [^abe] | anything but | alt("[^abe]") | abcde |
| [a-c] | range | alt("[a-c]") | abcde |

## ANCHORS

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

| regexp | matches | example | |
|---|---|---|---|
| ^a | start of string | anchor("^a") | aaa |
| a$ | end of string | anchor("a$") | aaa |

## LOOK AROUNDS

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

| regexp | matches | example | |
|---|---|---|---|
| a(?=c) | followed by | look("a(?=c)") | bacad |
| a(?!c) | not followed by | look("a(?!c)") | bacad |
| (?<=b)a | preceded by | look("(?<=b)a") | bacad |
| (?<!b)a | not preceded by | look("(?<!b)a") | bacad |

## QUANTIFIERS

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

| regexp | matches | example | |
|---|---|---|---|
| a? | zero or one | quant("a?") | .a.aa.aaa |
| a* | zero or more | quant("a*") | .a.aa.aaa |
| a+ | one or more | quant("a+") | .a.aa.aaa |
| a{n} | exactly **n** | quant("a{2}") | .a.aa.aaa |
| a{n, } | **n** or more | quant("a{2,}") | .a.aa.aaa |
| a{n, m} | between **n** and **m** | quant("a{2,4}") | .a.aa.aaa |

## GROUPS

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

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

| regexp | matches | example | |
|---|---|---|---|
| (ab\|d)e | 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) | regexp (to mean this) | matches (which matches this) | example (the result is the same as ref("abba")) | |
|---|---|---|---|---|
| \\1 | \1 (etc.) | first () group, etc. | ref("(a)(b)\\2\\1") | abbaab |

---

stringr