



## Workshop: Working with Strings in R

---

Simon Ress | Ruhr-Universität Bochum

October 18, 2021

# Content

## 1. Base R

The base R functions for dealing with strings

The base R functions for dealing regular expressions

## 2. Regular Expressions

## 3. Package: Stringr

Detect Matches

Subset Strings

Manage Lengths

Mutate Strings

## 4. Helpful sources

# Base R

---

# The base R functions for dealing with strings

- *substr()* / *substring()*: Extract or replace substrings in a character vector by indices.
- *strsplit()*: Split the elements of a character vector *x* into substrings according to a given regular expression.
- *paste()*: Concatenate *n* number of strings.
- *nchar()*: Returns a vector of the number of characters of *x*.

## substr() / substring()

### substr()/substring():

Extract or replace substrings in a character vector

```
num <- "12345678"
```

```
substr(num, 4, 5)
```

```
## [1] "45"
```

```
substring(num, 1:3, 7)
```

```
## [1] "1234567" "234567" "34567"
```

# strsplit()

## strsplit():

Split the elements of a character vector `x` into substrings according to a given character

```
str = "Splitting sentence into words"  
strsplit(str, " ")
```

```
## [[1]]
```

```
## [1] "Splitting" "sentence"  "into"      "words"
```

# paste()

**paste():**

Concatenate n number of strings

```
paste("Count number", "of characters")
```

```
## [1] "Count number of characters"
```

# nchar()

**nchar():**

Returns a vector of the number of characters of x

```
nchar("Count number of characters")
```

```
## [1] 26
```



# The base R functions for dealing regular expressions

- *grep()* / *grepl()*: Search for matches of a regular expression/pattern in a character vector and return the indices/a logical vector.
- *regexpr()* / *gregexpr()*: Search a character vector for regular expression matches and return the indices of the string where the match begins and the length of the match.
- *sub()* / *gsub()*: Search the first/all character vector/s for regular expression matches and replace that match with another string.

## grep() / grepl()

### grep():

Index of vector which matches regex

```
grep("b+", c("abc", "bda", "cca a", "abd"))
```

```
## [1] 1 2 4
```

### grepl():

Logical if vector matches regex

```
grepl("b+", c("abc", "bda", "cca a", "abd"))
```

```
## [1] TRUE TRUE FALSE TRUE
```

# regexpr()

## regexpr():

Search a character vector for regular expression matches and return the indices of the string where the match begins and the length of the match

```
str = "Line 129: 0 that this too too solid flesh would melt  
regexpr("1",str)
```

```
## [1] 6  
## attr(,"match.length")  
## [1] 1  
## attr(,"index.type")  
## [1] "chars"  
## attr(,"useBytes")  
## [1] TRUE
```

## sub() / gsub()

### sub():

Search **\*\*first\*** match of a regular expression and replace it

```
x <- "<dd>Found on January 1, 2007</dd>"
```

```
sub("<dd>[F|f]ound on |</dd>", "", x)
```

```
## [1] "January 1, 2007</dd>"
```

### gsub:

Search **\*\*all\*\*** matches of a regular expression and replace it

```
x <- "<dd>Found on January 1, 2007</dd>"
```

```
gsub("<dd>[F|f]ound on |</dd>", "", x)
```

```
## [1] "January 1, 2007"
```

# Regular Expressions

---

# Match Characters

## MATCH CHARACTERS

see <- function(rx) str\_view\_all("abc ABC 123\t.!?\\()\n", rx)

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

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

# Alternates

## ALTERNATES

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

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

# Anchors

## ANCHORS



regex

`^a`

`a$`

matches

start of string

end of string

example

`anchor("^a")`

`anchor("a$")`

aaa

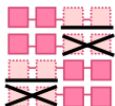
aaa

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



# Look Arouds

## LOOK AROUNDS



regex

`a(?=c)`

`a(?!c)`

`(?<=b)a`

`(?<!b)a`

matches

followed by

not followed by

preceded by

not preceded by

example

`look("a(?=c)")`

`look("a(?!c)")`

`look("(?<=b)a")`

`look("(?<!b)a")`

bacad

bacad

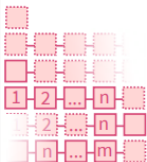
bacad

bacad

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

# Quantifiers

## QUANTIFIERS



regex

**a?**

**a\***

**a+**

**a{n}**

**a{n,}**

**a{n,m}**

matches

zero or one

zero or more

one or more

exactly **n**

**n** or more

between **n** and **m**

example

`quant("a?")`

`quant("a*")`

`quant("a+")`

`quant("a{2}")`

`quant("a{2,}")`

`quant("a{2,4}")`

`.a.aa.aaa`

`.a.aa.aaa`

`.a.aa.aaa`

`.a.aa.aaa`

`.a.aa.aaa`

`.a.aa.aaa`

# Package: Stringr

---

The stringr package provides a series of functions implementing much of the regular expression functionality in R but with a more consistent and rationalized interface.

# Detect Matches



→  
TRUE  
TRUE  
FALSE  
TRUE

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



→  
TRUE  
TRUE  
FALSE  
TRUE

**str\_starts**(string, **pattern**, negate = FALSE)  
Detect the presence of a pattern match at the beginning of a string. Also **str\_ends()**.  
`str_starts(fruit, "a")`



→  
1  
2  
4  
4

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



start end  
→  
2 4  
4 7  
NA NA  
3 4

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



→  
0  
3  
1  
2

**str\_count**(string, **pattern**) Count the number of matches in a string. `str_count(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, negate = FALSE)**  
Return only the strings that contain a pattern match. `str_subset(fruit, "p")`



**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(sentences, 6)`



**str\_trim(string, side = c("both", "left", "right"))** Trim whitespace from the start and/or end of a string. `str_trim(str_pad(fruit, 17))`

**str\_squish(string)** Trim whitespace from each end and collapse multiple spaces into single spaces. `str_squish(str_pad(fruit, 17, "both"))`

# 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. Also **str\_remove()**.

```
str_replace(fruit, "p", "-")
```



**str\_replace\_all()**(string, **pattern**, replacement)  
Replace all matched patterns in each string. Also **str\_remove\_all()**.

```
str_replace_all(fruit, "p", "-")
```

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. Also **str\_to\_sentence()**.  
**str\_to\_title()**(sentences)



## Extra: Join and Split



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



**str\_flatten(string, collapse = "")** Combines into a single string, separated by collapse.  
`str_flatten(fruit, ", ")`



**str\_dup(string, times)** Repeat strings times times. Also **str\_unique()** to remove duplicates.  
`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 and **str\_split\_n()** to return the nth substring.  
`str_split_fixed(sentences, " ", n=3)`



**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")`

## Extra: Order Strings



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



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

## Extra: Helpers

```
appl<e>  
banana  
p<e>ar
```



TRUE  
TRUE  
FALSE  
TRUE

This is a long sentence.

↓  
This is a long  
sentence.

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

**str\_view\_all**(string, **pattern**, match = NA)  
View HTML rendering of all regex matches.  
Also **str\_view()** to see only the first match.  
`str_view_all(sentences, "[aeiou]")`

**str\_equal**(x, y, locale = "en", ignore\_case = FALSE, ...)<sup>1</sup> Determine if two strings are equivalent. `str_equal(c("a", "b"), c("a", "c"))`

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

## Helpful sources

---

# Helpful sources

- [Stringr: Overview](#)
- [Stringr: Introduction](#)
- [Stringr: Cheatsheet](#)
- [Stringr: Reference manual](#)
- [Base R String-functions vs Stringr](#)
- [Working with strings in R](#)
- [Regular expressions](#)
- [Primary R functions for dealing with regular expressions](#)

All graphics are taken from String manipulation with stringr  
Cheatsheet