

Manipulating Strings: Brief Summary

Base R

Syntax (String = character string)	Description
<code>nchar(String)</code>	number of characters
<code>substr(String, start.m, end.n)</code>	substring from m to n
<code>strsplit(String, "w")</code>	split expression at (character) w
<code>sub("old", "new", String)</code>	Replace single instance of "old" with "new"
<code>gsub("old", "new", String)</code>	Replace instances of "old" with "new"

library(stringr)

Command	Description
<code>str_length(String)</code>	number of characters (similar to <code>nchar</code> .)
<code>str_sub(String, start.m, end.n)</code>	substring from m to n similar to <code>substr</code> .
<code>str_split(String, "w")</code>	split expression at (character) w
<code>str_detect(String, "pattern")</code>	detect presence/absence of pattern (return T/F)
<code>str_locate(String, "pattern")</code>	locate first position of pattern, return a matrix with columns start and end.
<code>str_locate_all(String, "pattern")</code>	locates all matches of pattern in string
<code>str_extract(String, "pattern")</code>	extracts text corresponding to first match
<code>str_extract_all(String, "pattern")</code>	extracts all matches and returns a list
<code>str_replace(String, "old", "new")</code>	Replace single instance of "old" with "new" (similar to <code>sub</code>).
<code>str_replace_all(String, "old", "new")</code>	Replace instances of "old" with "new" (similar to <code>gsub</code> .)

Regular Expressions

Metacharacters:

`. ^ \ $? * + [] () { } |`

If you actually want to use these symbols literally, then **in R**, you precede them by a double slash (or surround them with brackets).

```
str_detect(string, "\\.") #to match a period.  
str_detect(string, "[.]") #same
```

Note that outside of R , a single backslash works: `\$` to match the dollar sign.

Character sets

To match one of several options, use the brackets:

```
str_detect(string, "Ch[aio]mp")
```

matches Champ, Chimp, Chomp, Champion, Chimps, ...

Anchors

Use `^` to indicate a pattern that starts a string, `$` the end.

```
str_detect(String, "^Hello\\s") #string starts with hello (followed by space)
str_detect(String, "goodbye\\.?$") #string ends with goodbye (followed by period)
```

Syntax	Description
<code>\\d</code>	Digit, 0,1,2 ... 9
<code>\\D</code>	Not Digit
<code>\\s</code>	Space
<code>\\S</code>	Not Space
<code>\\w</code>	Word
<code>\\W</code>	Not Word
<code>\\t</code>	Tab
<code>\\n</code>	New line
<code>^</code>	Beginning of the string
<code>\$</code>	End of the string
<code>\\</code>	Escape special characters, e.g. <code>\\is "</code> , <code>\\+ is "+</code>
<code> </code>	Alternation match. e.g. <code>/(e d)n/</code> matches "en" and "dn"
<code>.</code>	Any character, except <code>\\n</code> or line terminator
<code>[ab]</code>	a or b
<code>[^ab]</code>	Any character except a and b
<code>[3-5]</code>	Range: for example, 3, 4, 5
<code>[C-E]</code>	Range: for example C, D, E
<code>[d-g]</code>	Range: for example, d, e, f or g
<code>[A-z]</code>	Range: Uppercase and lowercase a to z letters
<code>s+</code>	s at least one time
<code>s*</code>	s zero or more times
<code>s?</code>	s zero or 1 time
<code>s{n}</code>	s occurs n times in sequence
<code>s{n1,n2}</code>	s occurs between n1 and n2 times in sequence
<code>s{n1,n2}?</code>	non greedy match, see above example
<code>s{n,}</code>	s occurs n or more times
<code>()</code>	grouping
<code>.?*A</code>	match up to the first occurrence of A.

Some other useful patterns:

```
str_replace(String, "patternA(patternB)patternC", "\\1")
```

match above pattern and replace with patternB.