

Regular Expressions



Getting and Cleaning Data

A string is a sequence of characters, letters, numbers or symbols.

```
stringA <- "This sentence is a string."  
objectA <- c( "This sentence is a string.",  
              "Here's another string",  
              "Here's a third string" )
```



Regular expressions ("regexps")
are used to describe patterns
within strings.



```
first_string <- "The quick brown fox jumps over the lazy dog."  
second_string <- "There are 8 planets in our solar system:  
  Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, and Neptune.  
  There used to be a 9th planet: Pluto."
```

`str_view_all(string, pattern)`



```
str_view_all(second_string, "e")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter,
Saturn, Uranus, and Neptune. There used
to be a 9th planet: Pluto.



```
str_view_all(second_string, "8")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter,
Saturn, Uranus, and Neptune. There used to
be a 9th planet: Pluto.



```
str_view_all(second_string, "2")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter,
Saturn, Uranus, and Neptune. There used to
be a 9th planet: Pluto.



```
# Search for any letter
```

```
str_view_all(first_string, "[[:alpha:]]")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(second_string, "[[:alpha:]]")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter,
Saturn, Uranus, and Neptune. There used to
be a 9th planet: Pluto.




```
# Search our string for upper and lower case characters or any number.  
str_view_all(first_string, "[[:upper:]]")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(first_string, "[[:lower:]]")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(second_string, "[[:digit:]]")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter, Saturn,
Uranus, and Neptune. There used to be a 9th
planet: Pluto.



```
# Search our string for spaces and/or tabs  
str_view_all(first_string, "[[:space:]]")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(second_string, "[[:space:]]")
```

There are 8 planets in our solar system: Mercury,
Venus, Earth, Mars, Jupiter, Saturn, Uranus, and
Neptune. There used to be a 9th planet: Pluto.

```
str_view_all(first_string, "[[:blank:]]")
```

The quick brown fox jumps over the lazy dog.



```
# Search our string for punctuation  
str_view_all(first_string, "[[:punct:]]")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(second_string, "[[:punct:]]")
```

There are 8 planets in our solar system: Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, and Neptune. There used to be a 9th planet: Pluto.



```
# Search our string for any character but newline  
str_view_all(first_string, ".")
```

The quick brown fox jumps over the lazy dog.

```
str_view_all(second_string, ".")
```

There are 8 planets in our solar system:
Mercury, Venus, Earth, Mars, Jupiter,
Saturn, Uranus, and Neptune. There used to
be a 9th planet: Pluto.

```
# Search our string for any vowel  
str_view_all(first_string, "[aeiou]")
```

The quick brown fox jumps over the lazy dog.




```
# Search our string for all letters between a and m  
str_view_all(second_string, "[a-m]")
```

There are 8 planets in our solar system: Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, and Neptune. There used to be a 9th planet: Pluto.



```
# Search our string for NOT vowels  
str_view_all(first_string, "[^aeiou]")
```

The quick brown fox jumps over the lazy dog.



```
# Search the first string for NOT vowels OR spaces OR punctuation  
str_view_all(first_string, "[^aeiou | [[:space:]] | [[:punct:]] ]")
```

The quick brown fox jumps over the lazy dog.




```
names <-c("Keisha McDonald", "Mohammed Smith",  
          "Jane Doe", "Mathieu Person")
```

```
# Search for strings beginning with M.  
str_view_all(names, "^M")
```

Keisha McDonald

Mohammed Smith

Jane Doe

Mathieu Person



```
# Search for periods  
str_view_all(second_string, "\\.")
```

There are 8 planets in our solar system: Mercury, Venus, Earth, Mars, Jupiter, Saturn, Uranus, and Neptune. There used to be a 9th planet: Pluto.



```
# Create a new string with some special characters in it  
crazy_string <- "This string is cRaZy!?! :)"
```

```
# Search it for !, ?, and )  
str_view_all(crazy_string, "\\!")
```

This string is cRaZy!?! :)

```
str_view_all(crazy_string, "\\?")
```

This string is cRaZy!?! :)

```
str_view_all(crazy_string, "\\)")
```

This string is cRaZy!?! :)



Quantifiers

- $? : 0 \text{ or } 1$
- $+ : 1 \text{ or more}$
- $* : 0 \text{ or more}$

- $\{n\} : \text{exactly } n \text{ times}$
- $\{n,\} : n \text{ or more times}$
- $\{n,m\} : \text{between } n \text{ and } m \text{ times}$



```
# Identify any time the letter m shows up one or more times  
str_view_all(names, "m+")
```

Keisha McDonald

Mohammed Smith

Jane Doe

Mathieu Person



```
# Identify single m's  
str_view_all(names, "m{1}")
```

Keisha McDonald

Mohammed Smith

Jane Doe

Mathieu Person



```
# Identify where m shows up exactly two times in a row  
str_view_all(names, "m{2}")
```

Keisha McDonald

Mohammed Smith

Jane Doe

Mathieu Person



```
# Identify any time 'mm' shows up one or more times  
str_view_all(names, "mm+")
```

Keisha McDonald

Mohammed Smith

Jane Doe

Mathieu Person



Summarizing: Regular Expressions



Getting and Cleaning Data

Basic Regular Expressions in R

Cheat Sheet

Character Classes

| | |
|---|---|
| <code>[[digit:]]</code> or <code>\d</code> | Digits; [0-9] |
| <code>\D</code> | Non-digits; [^0-9] |
| <code>[[lower:]]</code> | Lower-case letters; [a-z] |
| <code>[[upper:]]</code> | Upper-case letters; [A-Z] |
| <code>[[alpha:]]</code> | Alphabetic characters; [A-z] |
| <code>[[alnum:]]</code> | Alphanumeric characters [A-z0-9] |
| <code>\w</code> | Word characters; [A-z0-9_] |
| <code>\W</code> | Non-word characters |
| <code>[[xdigit:]]</code> or <code>\X</code> | Hexadec. digits; [0-9A-Fa-f] |
| <code>[[blank:]]</code> | Space and tab |
| <code>[[space:]]</code> or <code>\s</code> | Space, tab, vertical tab, newline, form feed, carriage return |
| <code>\S</code> | Not space; [^[:space:]] |
| <code>[[punct:]]</code> | Punctuation characters; [!\"#\$%&'()*+,-./:;<=>?@[\^_`{ }~ |
| <code>[[graph:]]</code> | Graphical characters; [[:alnum:]][[:punct:]] |
| <code>[[print:]]</code> | Printable characters; [[:alnum:]][[:punct:]]\s |
| <code>[[cntrl:]]</code> or <code>\c</code> | Control characters; \n, \r etc. |

Special Metacharacters

| | |
|-----------------|-----------------|
| <code>\n</code> | New line |
| <code>\r</code> | Carriage return |
| <code>\t</code> | Tab |
| <code>\v</code> | Vertical tab |
| <code>\f</code> | Form feed |

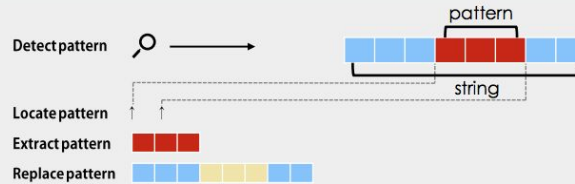
Lookaheads and Conditionals*

| | |
|----------------------------|--|
| <code>(?=)</code> | Lookahead (requires PERL = TRUE), e.g. <code>(?=xy)</code> : position followed by 'xy' |
| <code>(?!)</code> | Negative lookahead (PERL = TRUE); position NOT followed by pattern |
| <code>(?<=)</code> | Lookbehind (PERL = TRUE), e.g. <code>(?<=xy)</code> : position following 'xy' |
| <code>(?<!)</code> | Negative lookbehind (PERL = TRUE); position NOT following pattern |
| <code>?(!)then</code> | If-then-condition (PERL = TRUE); use lookaheads, optional char. etc in if-clause |
| <code>?(!)then else</code> | If-then-else-condition (PERL = TRUE) |

*see, e.g. <http://www.regular-expressions.info/lookaround.html>
<http://www.regular-expressions.info/conditional.html>

CC BY Ian Kopacka • ian.kopacka@ages.at

Functions for Pattern Matching



```
> string <- c("Hiphopotamus", "Rhyenoceros", "time for bottomless lyrics")
> pattern <- "t.m"
```

Detect Patterns

```
grep(pattern, string)
[1] 1 3

grep(pattern, string, value = TRUE)
[1] "Hiphopotamus"
[2] "time for bottomless lyrics"

grepl(pattern, string)
[1] TRUE FALSE TRUE

stringr::str_detect(string, pattern)
[1] TRUE FALSE TRUE
```

Split a String using a Pattern

`strsplit(string, pattern)` or `stringr::str_split(string, pattern)`

Locate Patterns

```
regexpr(pattern, string)
find starting position and length of first match

gregexpr(pattern, string)
find starting position and length of all matches

stringr::str_locate(string, pattern)
find starting and end position of first match

stringr::str_locate_all(string, pattern)
find starting and end position of all matches
```

Extract Patterns

```
regmatches(string, regexpr(pattern, string))
extract first match [1] "tam" "tim"

regmatches(string, gregexpr(pattern, string))
extract all matches, outputs a list
[[1]] "tam" [[2]] character(0) [[3]] "tim" "tom"

stringr::str_extract(string, pattern)
extract first match [1] "tam" NA "tim"

stringr::str_extract_all(string, pattern)
extract all matches, outputs a list

stringr::str_match(string, pattern)
extract first match + individual character groups

stringr::str_match_all(string, pattern)
extract all matches + individual character groups
```

Replace Patterns

```
sub(pattern, replacement, string)
replace first match

gsub(pattern, replacement, string)
replace all matches

stringr::str_replace(string, pattern, replacement)
replace first match

stringr::str_replace_all(string, pattern, replacement)
replace all matches
```

Character Classes and Groups

| | |
|---------------------|--|
| <code>.</code> | Any character except \n |
| <code> </code> | Or, e.g. <code>(a b)</code> |
| <code>[...]</code> | List permitted characters, e.g. <code>[abc]</code> |
| <code>[a-z]</code> | Specify character ranges |
| <code>[^...]</code> | List excluded characters |
| <code>(...)</code> | Grouping, enables back referencing using <code>\N</code> where N is an integer |

General Modes

By default R uses *extended* regular expressions. You can switch to *PCRE regular expressions* using `PERL = TRUE` for base or by wrapping patterns with `perl()` for stringr.

All functions can be used with literal searches using `fixed = TRUE` for base or by wrapping patterns with `fixed()` for stringr.

All base functions can be made case insensitive by specifying `ignore.case = TRUE`.

Anchors

| | |
|--------------------|---------------------------------------|
| <code>^</code> | Start of the string |
| <code>\$</code> | End of the string |
| <code>\b</code> | Empty string at either edge of a word |
| <code>\B</code> | NOT the edge of a word |
| <code>\<</code> | Beginning of a word |
| <code>\></code> | End of a word |

Escaping Characters

Metacharacters (`.`, `*`, `+` etc.) can be used as literal characters by escaping them. Characters can be escaped using `\\` or by enclosing them in `\\Q...\\E`.

Case Conversions

Regular expressions can be made case insensitive using `(?i)`. In backreferences, the strings can be converted to lower or upper case using `\\L` or `\\U` (e.g. `\\L\\l\\1`). This requires `PERL = TRUE`.

Quantifiers

| | |
|--------------------|---|
| <code>*</code> | Matches at least 0 times |
| <code>+</code> | Matches at least 1 time |
| <code>?</code> | Matches at most 1 time; optional string |
| <code>{n}</code> | Matches exactly n times |
| <code>{n,}</code> | Matches at least n times |
| <code>{n,m}</code> | Matches between n and m times |

Greedy Matching

By default the asterisk `*` is greedy, i.e. it always matches the longest possible string. It can be used in lazy mode by adding `?`, i.e. `*?`.

Greedy mode can be turned off using `(?U)`. This switches the syntax, so that `(?U)a*` is lazy and `(?U)a*?` is greedy.

Note

Regular expressions can conveniently be created using e.g. the packages `rex` or `rebus`.

Updated: 10/18

