

# 14 - Strings

ST 597 | Spring 2016  
University of Alabama

*14-strings.pdf*

## Contents

<b>1</b>	<b>Strings and Characters</b>	<b>2</b>
1.1	Text as Data . . . . .	2
<b>2</b>	<b>stringr R package</b>	<b>3</b>
2.1	Combining and Collapsing Strings . . . . .	3
2.2	String Subsetting . . . . .	5
2.3	String Length . . . . .	5
2.4	String Formatting . . . . .	7
<b>3</b>	<b>Regular Expressions: Finding Patterns</b>	<b>8</b>
3.1	Literal String Matching . . . . .	9
3.2	Multiple Choice . . . . .	10
3.3	Shortcuts . . . . .	10
3.4	Quantifiers . . . . .	11
3.5	Grouping . . . . .	12
3.6	Greedy Matching . . . . .	12
3.7	Anchors . . . . .	12
<b>4</b>	<b>String Manipulation</b>	<b>13</b>
4.1	Remove text . . . . .	13
4.2	Finding all words . . . . .	13
<b>5</b>	<b>Pattern Detection</b>	<b>14</b>
5.1	Word Count . . . . .	14
5.2	Baby Names . . . . .	15
<b>6</b>	<b>Examples</b>	<b>16</b>
6.1	Baseball 3,000 Hit Club . . . . .	16

---

## Required Packages and Data

```
library(tidyverse)
library(stringr)
library(rvest)
library(babynames)
#-
# library(devtools) # install.packages("devtools") if necessary
# devtools::install_github("hadley/stringr")
```

## 1 Strings and Characters

- A **character** is a single symbol
  - one byte per character
  - letter, number, punctuation, whitespace
  - special characters (`\n`, `\r`, `\t`)
  - encoding
- A **string** is a collection of 0 or more characters
  - "Bacon and Eggs" is a string of letters and whitespaces
  - "" is an empty string
  - "www.ua.edu"
  - "yourname@crimson.ua.edu"

### 1.1 Text as Data

Think of text (strings) as data, just like you do for numbers. In R, strings are considered *character* data.

```
x = "Bacon and Eggs"
class(x)
#> [1] "character"
```

In R, the string `x` is a one element vector. Here is a vector of three strings:

```
y = c("Bacon and Eggs", "French Toast", "coffee")
length(y)
#> [1] 3
y[2:3]
#> [1] "French Toast" "coffee"
```

Usually, single or double quotes will suffice.

```
single = 'The fear of the Lord is the beginning of wisdom'
double = "The fear of the Lord is the beginning of wisdom"
identical(single, double)
#> [1] TRUE
```

An exception is with quotations or apostrophes. The `\` is an *escape* symbol that forces the next character to be taken literally.

```
quotes = "She said, \"there's is no time like the present\"" # escape \
quotes2 = 'She said, "there's is no time like the present"' # mix quotes
```

```
quotes
#> [1] "She said, \"there's is no time like the present\""
writeLines(quotes)           # writeLines() displays in readable form
#> She said, "there's is no time like the present"
identical(quotes, quotes2)
#> [1] TRUE
```

The `\` can also be used to indicate special things like tab (`\t`)

```
writeLines("there's no \time like the present")
#> there's no   ime like the present
```

## 2 stringr R package

<http://cran.r-project.org/web/packages/stringr/stringr.pdf>

```
library(stringr)
```

From the [Intro to stringr vignette](#):

Strings are not glamorous, high-profile components of R, but they do play a big role in many data cleaning and preparations tasks. R provides a solid set of string operations, but because they have grown organically over time, they can be inconsistent and a little hard to learn. Additionally, they lag behind the string operations in other programming languages, so that some things that are easy to do in languages like Ruby or Python are rather hard to do in R. The **stringr** package aims to remedy these problems by providing a clean, modern interface to common string operations.

More concretely, stringr:

- Simplifies string operations by eliminating options that you don't need 95% of the time (the other 5% of the time you can use functions from base R or [stringi](#)).
- Uses consistent function names and arguments.
- Produces outputs that can easily be used as inputs. This includes ensuring that missing inputs result in missing outputs, and zero length inputs result in zero length outputs. It also processes factors and character vectors in the same way.
- Completes R's string handling functions with useful functions from other programming languages.

To meet these goals, stringr provides two basic families of functions:

- basic string operations, and
- pattern matching functions which use regular expressions to detect, locate, match, replace, extract, and split strings.

As of version 1.0, stringr is a thin wrapper around [stringi](#), which implements all the functions in stringr with efficient C code based on the [ICU library](#). Compared to stringi, stringr is considerably simpler: it provides fewer options and fewer functions. This is great when you're getting started learning string functions, and if you do need more of stringi's power, you should find the interface similar.

### 2.1 Combining and Collapsing Strings

function	description
<code>str_c()</code>	Join multiple strings into a single string
<code>paste()</code>	Concatenate vectors after converting to character
<code>str_dup()</code>	Duplicate and concatenate strings within a character vector

The `stringr` function `str_c()` is basically the base R function `paste0()`. We will use `str_c()` just to stay consistent.

### 2.1.1 Combining vectors into strings

```
str_c("A", "B", "C")
#> [1] "ABC"

#- change the sep= argument to a space
str_c("A", "B", "C", sep=" ")
#> [1] "A B C"

#- change the sep= argument to a :
str_c("A", "B", "C", sep=":")
#> [1] "A:B:C"
```

Vectorized (remember vector recycling)

```
str_c("X", 1:5)
#> [1] "X1" "X2" "X3" "X4" "X5"
ABC = c("A", "B", "C")
str_c(ABC, 1:3)
#> [1] "A1" "B2" "C3"
str_c(ABC, 1:3, sep='-')
#> [1] "A-1" "B-2" "C-3"
str_c(LETTERS[1:6], 1:3, sep='-')
#> [1] "A-1" "B-2" "C-3" "D-1" "E-2" "F-3"

str_c("---", "Iteration: ", 1:10, "---" )
#> [1] "---Iteration: 1---" "---Iteration: 2---" "---Iteration: 3---"
#> [4] "---Iteration: 4---" "---Iteration: 5---" "---Iteration: 6---"
#> [7] "---Iteration: 7---" "---Iteration: 8---" "---Iteration: 9---"
#> [10] "---Iteration: 10---"
```

There is also the `str_dup()` function (compare to `rep()`)

```
str_dup(ABC, times = 2)
#> [1] "AA" "BB" "CC"
str_dup(ABC, times = 1:3)
#> [1] "A" "BB" "CCC"

rep(ABC, times=1:3)
#> [1] "A" "B" "B" "C" "C" "C"
rep(ABC, each=2)
#> [1] "A" "A" "B" "B" "C" "C"
```

## 2.1.2 Make one long string

The `collapse=` argument of `str_c()` (and `paste()`),

```
str_c(ABC)           # three element vector
#> [1] "A" "B" "C"
str_c(ABC, collapse="") # one element vector
#> [1] "ABC"
str_c(ABC, collapse=", ") # one element vector
#> [1] "A, B, C"
```

## 2.2 String Subsetting

function	description
<code>str_sub()</code>	Extract and replace substrings from a character vector

We have already used the function `str_sub()` to subset strings by position.

```
data(state) # access to state.name data
str_sub(state.name, start=1, end=4) # 1st 4 letters
#> [1] "Alab" "Alas" "Ariz" "Arka" "Cali" "Colo" "Conn" "Dela" "Flor" "Geor" "Hawa" "Idah"
#> [13] "Illi" "Indi" "Iowa" "Kans" "Kent" "Loui" "Main" "Mary" "Mass" "Mich" "Minn" "Miss"
#> [25] "Miss" "Mont" "Nebr" "Neva" "New " "New " "New " "New " "Nort" "Nort" "Ohio" "Okla"
#> [37] "Oreg" "Penn" "Rhod" "Sout" "Sout" "Tenn" "Texa" "Utah" "Verm" "Virg" "Wash" "West"
#> [49] "Wisc" "Wyom"
str_sub(state.name, start=-4, end=-1) # last 4 letters
#> [1] "bama" "aska" "zona" "nsas" "rnia" "rado" "icut" "ware" "rida" "rgia" "waii" "daho"
#> [13] "nois" "iana" "Iowa" "nsas" "ucky" "iana" "aine" "land" "etts" "igan" "sota" "ippi"
#> [25] "ouri" "tana" "aska" "vada" "hire" "rsey" "xico" "York" "lina" "kota" "Ohio" "homa"
#> [37] "egon" "ania" "land" "lina" "kota" "ssee" "exas" "Utah" "mont" "inia" "gton" "inia"
#> [49] "nsin" "ming"
```

## 2.3 String Length

function	description
<code>str_length()</code>	The length of a string

The function `str_length()` returns the number of characters (or *code points*) in a string:

```
tibble(state=state.name, letters = str_length(state.name))
#> # A tibble: 50 × 2
#>       state letters
#>   <chr>   <int>
#> 1 Alabama     7
#> 2 Alaska      6
#> 3 Arizona     7
#> 4 Arkansas    8
```

```
#> 5    California      10
#> 6    Colorado       8
#> 7    Connecticut    11
#> 8    Delaware       8
#> 9    Florida        7
#> 10   Georgia        7
#> # ... with 40 more rows
```

### 2.3.1 Characters in Colossians

We will further illustrate with a text version of the Book of Colossians from the bible. The text can be found here <https://raw.githubusercontent.com/mdporter/ST597/master/data/colossians.txt>.

```
url = "https://raw.githubusercontent.com/mdporter/ST597/master/data/colossians.txt"
lines = read_lines(url)      # vector: each element a string
single = read_file(url)      # single element: string of entire file
```

Notice that the vector `lines` has one element per line in the document, while `single` is one long string.

Let's see how many characters are in the document

```
str_length(lines)           # notice there are empty lines
#> [1]  91  0 116  0  93  0 108  0 128  0 197  0 129  0  63  0 184  0 156  0 1
#> [22]  0 117  0 116  0  67  0  83  0 201  0  80  0 151  0  71  0 137  0  86
#> [43] 145  0 227  0 179  0 148  0  90  0 163  0 134  0  98  0 145  0 191  0
#> [64]  0  86  0 149  0  79  0 122  0 193  0  67  0  90  0 156  0 169  0 167
#> [85] 136  0 113  0 145  0  93  0 169  0 174  0 154  0  58  0 113  0 195  0 1
#> [106]  0  87  0  76  0  96  0 148  0  61  0  73  0 114  0 101  0 112  0 143
#> [127] 133  0 162  0  99  0 125  0 197  0 142  0  73  0  72  0  84  0  83  0 1
#> [148]  0  79  0 123  0 105
sum(str_length(lines))      # total number of characters
#> [1] 9527
str_length(single)         # number of characters in long string
#> [1] 9679
```

Why do the number of characters not match?

Because `single` has new line (`\n`) and carriage return (`\r`) symbols.

```
lines[1:3]
#> [1] "[Col 1:1 ESV] Paul, an apostle of Christ Jesus by the will of God, and Timothy our
#> [2] ""
#> [3] "[Col 1:2 ESV] To the saints and faithful brothers in Christ at Colossae: Grace to y
str_sub(single, start=1, end=211)
#> [1] "[Col 1:1 ESV] Paul, an apostle of Christ Jesus by the will of God, and Timothy our
```

We can remove these with the `str_replace_all()` function.

```
single2 = str_replace_all(string=single, pattern="\n|\r", replacement = "" )
str_length(single2)
#> [1] 9527
sum(str_length(lines))
#> [1] 9527
```

```
str_sub(single2, 1, 211)
#> [1] "[Col 1:1 ESV] Paul, an apostle of Christ Jesus by the will of God, and Timothy our
```

The `|` symbol in `pattern="\n|\r"` means we want to remove all new line (`\n`) and carriage return (`\r`).

Side note: we can collapse the line-by-line vector `lines` to one long string.

```
lines2 = str_c(lines, collapse="")
identical(lines2, single2)
#> [1] TRUE
```

## 2.4 String Formatting

function	description
<code>str_wrap()</code>	Wrap strings into nicely formatted paragraphs
<code>str_trunc()</code>	Truncate a character string
<code>str_trim()</code>	Trim whitespace from start and end of string
<code>str_pad()</code>	Pad a string
<code>str_to_lower()</code>	Convert string to all lowercase: see <code>tolower()</code>
<code>str_to_upper()</code>	Convert string to all uppercase: see <code>toupper()</code>
<code>str_to_title()</code>	Convert string to all title case
<code>str_conv()</code>	Specify the encoding of a string

To help see long strings, the `str_wrap()` and `str_trunc()` are helpful

```
writeLines(lines[1]) # writeLines() displays in readable form
#> [Col 1:1 ESV] Paul, an apostle of Christ Jesus by the will of God, and Timothy our brot
writeLines(str_wrap(lines[1], width=50)) # adds newlines '\n'
#> [Col 1:1 ESV] Paul, an apostle of Christ Jesus by
#> the will of God, and Timothy our brother,
writeLines(str_wrap(lines[1], width=1)) # special treatment for width<=1
#> [Col
#> 1:1
#> ESV]
#> Paul,
#> an
#> apostle
#> of
#> Christ
#> Jesus
#> by
#> the
#> will
#> of
#> God,
#> and
#> Timothy
#> our
#> brother,
```

```
str_trunc(lines[1], width = 30)      # reduces string to only `width` characters
#> [1] "[Col 1:1 ESV] Paul, an apos..."
```

The functions `str_trim()` and `str_pad()` help with whitespace and length

```
str_trim("  String with trailing and leading white space\t")
#> [1] "String with trailing and leading white space"
str_trim("\n\nString with trailing and leading white space\n\n")
#> [1] "String with trailing and leading white space"

x = c("A", "BB", "CCC", "DDDD")
(x.pad = str_pad(x, width=4, side="right")) # `width` is minimum width
#> [1] "A   " "BB  " "CCC " "DDDD"
str_length(x.pad)
#> [1] 4 4 4 4
str_pad(x, width=4, side="right", pad="-")
#> [1] "A---" "BB--" "CCC-" "DDDD"
```

### 2.4.1 Convert the case of a string

For extracting words from text, we often want to remove the effects of case.

```
s1 = "This is not the same sentence"
s2 = "this is NOT the same sentence"

#- identical when all lowercase
str_to_lower(s1)
#> [1] "this is not the same sentence"
str_to_lower(s2)
#> [1] "this is not the same sentence"

#- identical when all uppercase
str_to_upper(s1)
#> [1] "THIS IS NOT THE SAME SENTENCE"
str_to_upper(s2)
#> [1] "THIS IS NOT THE SAME SENTENCE"

#- title case example
str_to_title(s1)
#> [1] "This Is Not The Same Sentence"
```

## 3 Regular Expressions: Finding Patterns

Regular expressions, regexps for short, are a very terse language that allow to describe patterns in strings. They take a little while to get your head around, but once you've got it you'll find them extremely useful. We have only used the most basic patterns so far:

- remove the and from 'Bacon and Eggs'



```
y = 'Bacon and Eggs'
str_replace(y, pattern = " and", replace="") # removed `whitespace + and`
#> [1] "Bacon Eggs"
```

- split the first and last names by finding a whitespace

```
student <- c("John Davis", "Angela Williams", "Bullwinkle Moose",
            "David Jones", "Janice Markhammer",
            "Cheryl Cushing", "Reuven Ytzrhak",
            "Greg Knox", "Joel England", "Mary Rayburn")
str_split_fixed(student, " ", n=2)
#>      [,1]      [,2]
#> [1,] "John"    "Davis"
#> [2,] "Angela"   "Williams"
#> [3,] "Bullwinkle" "Moose"
#> [4,] "David"    "Jones"
#> [5,] "Janice"   "Markhammer"
#> [6,] "Cheryl"   "Cushing"
#> [7,] "Reuven"   "Ytzrhak"
#> [8,] "Greg"     "Knox"
#> [9,] "Joel"     "England"
#> [10,] "Mary"    "Rayburn"
```

In these examples, the pattern was an actual symbol or set of symbols. But we need to be able to find patterns in a string (e.g., an uppercase letter followed by one or more lowercase, an optional three digit area code followed by 7 numbers, an email address). Regular expressions gives us a very flexible way to specify patterns.

We will not cover regular expressions in much detail. And it is a big topic. Here are some references to bookmark:

- <http://stackoverflow.com/questions/4736/learning-regular-expressions>
- [stringr vignette](#)
- [R regex help page](#)
- [Learning Regular Expressions](#)
- [Interactive Testing](#)
- [Build regular expression from an input string](#)
- [Regex Handout - Shalizi](#)
- [R help ?regex](#)

### 3.1 Literal String Matching

Consider the quote from Albert Einstein, “Not everything that can be counted counts, and not everything that counts can be counted.”

```
#- literal strings must be on one line, or R will insert a new line (\n) character
x = "Not everything that can be counted counts, and not everything that counts can be counted"
```

Note: if you want to write long string on multiple lines, use another function line `strwrap()`.

Regular expressions will find all consecutive characters that match the pattern. The simplest example is when

the pattern is a usual string:

```
str_view_all(x, 'counts')      # finds all `counts`
str_view_all(x, 'counted')     # finds all `counted`
str_view_all(x, 'count')       # finds the `count` in `counted` and `counts`
str_view_all(x, 'c')           # finds all `c`
str_view_all(x, 'a')           # finds all `a`
str_view_all(x, ' ')           # finds all whitespaces
str_view_all(x, 'can be')      # finds the `can be` string
str_view_all(x, 'not')         # finds all `not` with lowercase `n`
str_view_all(x, 'Not')         # finds all `Not` with uppercase `N`
str_view_all(x, 'not ever')    # does this find what you expect?
```

Notice that the pattern 'count' finds the sequence 'c', followed by a 'o', followed by a 'u', followed by a 'n', followed by a 't'.

## 3.2 Multiple Choice

You can specify multiple options (logical OR) with a vertical bar |. For example, to match 'Not' *or* 'not' we use the pattern 'Not|not'

```
str_view_all(x, pattern='Not|not')
```

### 3.2.1 Character Classes: Pick one

For *single characters*, you can indicate multiple options by enclosing in brackets ([ ]). For example, the previous example can also be obtained

```
str_view_all(x, pattern='[Nn]ot')  # matches 'N' or 'n', followed by 'ot'
str_view_all(x, pattern='[aeiou]')  # find all lowercase vowels
```

If the first character is the caret (^), then it finds any character **not** in the list. For example, to find all (lowercase) consonants

```
str_view_all(x, pattern='[^aeiou]') # find all lowercase consonants (not vowels)
```

## 3.3 Shortcuts

Some fruit and phone numbers:

```
fruit = c("apple", "banana", "pear", "pineapple")

strings = c(
  "apple",
  "219 733 8965",
  "329-293-8753",
  "Work: 579-499-7527; Home: 543.355.3679")
```

To help with common patterns, most regular expression engines recognize shortcuts. For example, ranges of characters are also recognized in character classes.

- `[a-d]` matches the letters 'a' or 'b' or 'c' or 'd'
- `[A-Z]` matches all uppercase letters
- `[0-5]` matches the numbers '0' or '1' or '2' ,,, or '5'

```
str_view_all(strings, pattern='[7-9]')
```

Other special shortcuts include:

- The period (`.`) matches any single character (including whitespace, but not newline)

```
str_view_all(fruit, pattern='a.') # finds 'a' followed by anything
str_view_all(x, pattern='t.')    # finds 't' followed by anything
```

- Type `?regex` to see the following shortcuts (there are more than these)
  - `[:digit:]` (or `\d`) is same as `[0-9]`
  - `[:lower:]` and `[:upper:]`, the upper and lower case letters. This is generally the same as `[a-z]` and `[A-Z]` except the former are locale independent.
  - `[:alpha:]` all letters of the alphabet. Union of `[:lower:]` and `[:upper:]`
  - `[:alnum:]` the alphanumeric characters. Union of `[:alpha:]` and `[:digit:]` and same as `[0-9A-Za-z]`
  - `[:punct:]` all the punctuation characters: `! " # $ % & ' ( ) * + , - . / : ; < = > ? @ [ \ ] ^ _ { | } ~ . ' , - . / : ;`
  - `[:space:]` (also `\s`) all the whitespace characters: tab, newline, space, carriage return

Some other common shortcuts are:

- `\w` stands for the word characters `[A-Za-z0-9_]` or `[:alnum:]` (alphanumeric *union* underscore)
- `\W` is anything but a word character
- `\d` any digit (same as `[:digit:]` and `[0-9]`)
- `\D` any non-digit
- `\s` any whitespace
- `\S` any non-whitespace
- `\w` any word character
- `\b` boundary
- We have already seen the new line (`\n`) and carriage return (`\r`)

But notice to use these in R, we need to use a double backslash because the first one is treated as an escape

```
str_view_all(strings, pattern='\\d') # finds digits
str_view_all(strings, pattern='\\s') # finds all whitespaces
str_view_all(strings, pattern='[\\s[:upper:]]') # whitespaces OR uppercase
```

### 3.4 Quantifiers

Quantifiers allow control over how many times a pattern is matched.

- `?` (zero or one time)
- `*` (zero or more times)
- `+` (one or more times)
- `{n}` (exactly n times)
- `{n, }` (at least n times)

- $\{, m\}$  (at most  $m$  times)
- $\{n, m\}$  (at least  $n$  times but no more than  $m$  times)

```
#- find the phone numbers
str_view_all(strings, pattern='\\d{3}-\\d{3}-\\d{4}')
```

```
str_view_all(strings, pattern='\\d{3}[- \\.]\\d{3}[- \\.]\\d{4}')
```

The first one finds the pattern XXX-XXX-XXXX where X is any digit. The second allows several separators.

### 3.5 Grouping

A quantifier modifies the character (or shortcut) to its immediate left. So the pattern `0abc+0` matches `'0abc0'`, `'0abcc0'`, `'0abccc0'`, etc. but not `'0abcabc0'`, `'0abcabcabc0'`, etc. To get this behavior just wrap in parentheses. E.g., `0(abc)+0` will find `abc` one or more times.

```
y = c("0abc0", "0abcc0", "0abccc0", "0abcabc0", "0abcabcabc0")

str_view_all(y, pattern="0abc+0")
str_view_all(y, pattern="0(abc)+0")
```

### 3.6 Greedy Matching

Regular expressions will match as much as it can. For example, say we want to find the `ggplot()` function in an R script.

```
R = "ggplot(cars, aes(reorder(type, enginesize)))+geom_bar()+facet_wrap(~origin) +
labs(x='car type')"
```

The pattern is `ggplot` followed by a `(`, the varying set of arguments, then ends in a `)`. We may want to try

```
str_view_all(R, pattern='ggplot\\((.*\\)\\)') # entire string matched
```

Notice that it matches to the last `)` in the entire string. This is technically correct, but not what we want.

To match the first occurrence, we can add an extra `?` to the quantifier. So for example,

```
str_view_all(R, pattern='ggplot\\((.*?\\)\\)') # finds first closing parenthesis
str_view_all(R, pattern='ggplot\\((.*?\\)+\\)') # finds first set of closing parentheses
```

The first example stops too early. We really want to match to the last of a set of closing parentheses.

### 3.7 Anchors

Regular expressions match any part of a string. So the pattern `apple` matches the `apple` in the fruit `'apple'` as well as the last part of `'pineapple'`

```
str_view_all(fruit, pattern='apple') # finds all 'apple' substrings
```

We can limit to finding all strings that start with `'apple'` by using the `^`

```
str_view_all(fruit, pattern='^apple') # finds strings that start with 'apple'
```

The special pattern `^` matches only at the beginning of your input and `$` to match only at the end. This allows *bookends* where you can specify what is at the beginning and end.

```
str_view_all(fruit, pattern='apple$') # finds strings that end with `apple`
```

## 4 String Manipulation

function	description
<code>str_replace()</code>	Replace matched patterns in a string
<code>str_replace_na()</code>	Turn NA into “NA”
<code>str_sub()</code>	Extract and replace substrings from a character vector
<code>str_split()</code>	Split up a string into pieces
<code>str_split_fixed()</code>	Split up a string into exactly <i>n</i> pieces

- Many of these functions have a version to find *all* matching patterns using `str_<function>_all`.

Recall the bible verses from the Book of Colossians

```
url = "https://raw.githubusercontent.com/mdporter/ST597/master/data/colossians.txt"
lines = read_lines(url) # vector: each element a string
single = read_file(url) # single element: string of entire file
```

Notice that the vector `lines` has one element per line in the document, while `single` is one long string.

### 4.1 Remove text

We have already seen most of these functions in action. Let's use `str_replace_all()` to remove the verse reference.

```
col_text = str_replace_all(single, pattern='\\[.+? ESV\\]', replacement='')
```

### 4.2 Finding all words

One useful task is to extract all words from a string. We can use the `str_split()` function to break a string up into components. The helper function `boundary(type = "word")` will split a string into word components, dropping the punctuation.

```
str_split(col_text, " ") [[1]] %>% head(30)
#> [1] "" "Paul," "an" "apostle" "of"
#> [6] "Christ" "Jesus" "by" "the" "will"
#> [11] "of" "God," "and" "Timothy" "our"
#> [16] "brother,\n\n" "To" "the" "saints" "and"
#> [21] "faithful" "brothers" "in" "Christ" "at"
#> [26] "Colossae:" "Grace" "to" "you" "and"
str_split(col_text, boundary("word")) [[1]] %>% head(30)
#> [1] "Paul" "an" "apostle" "of" "Christ" "Jesus" "by"
#> [8] "the" "will" "of" "God" "and" "Timothy" "our"
```

```
#> [15] "brother" "To" "the" "saints" "and" "faithful" "brothers"
#> [22] "in" "Christ" "at" "Colossae" "Grace" "to" "you"
#> [29] "and" "peace"
```

Notice that `str_split()` will return a length one *list* - and we need to extract the first component using `[[1]]` (which will return a vector of words).

## 5 Pattern Detection

function	description
<code>str_detect()</code>	Detect the presence or absence of a pattern in a string
<code>str_count()</code>	Count the number of matches in a string
<code>str_extract()</code>	Extract matching patterns from a string
<code>str_match()</code>	Extract matched groups from a string
<code>str_subset()</code>	Keep strings matching a pattern
<code>str_locate()</code>	Locate the position of patterns in a string

- Many of these functions have a version to find *all* matching patterns using `str_<function>_all`.

### 5.1 Word Count

The function `str_count()` counts the number of times the pattern is matched. Here we will extract the number of times the pattern `God` is found

```
str_count(single, 'God')
#> [1] 20
```

Notice that the *pattern* `God` would match `God's`, `Godspeed`, etc.

We can get a count of all words using the combination of `str_split()`, `tibble()` and `count()`

```
#- get the words from Book of Colossians
col_words = str_split(col_text, boundary("word"))[[1]] # vector of words

#- Get frequency of words
word_freq = tibble(word=col_words) %>% count(word, sort=TRUE)
```

But here we see that `God` is matched only 18 times (not 20). This is because our regular expression pattern matched multiple *words*.

```
filter(word_freq, str_detect(word, 'God'))
#> # A tibble: 2 × 2
#>   word      n
#>   <chr> <int>
#> 1  God     18
#> 2 God's     2
```

## 5.2 Baby Names

The function `str_subset()`, which is equivalent to `x[str_detect(x, pattern)]`, returns the strings with matching pattern.

Use the `babynames` data from the `babynames` package to find variations of your name.

```
# update package install.packages("babynames") to 2014
library(babynames)

#- get names from 2014 (only most common names included)
baby2014 = filter(babynames, year==2014)

library(stringr)

#- return matching strings
str_subset(baby2014$name, "Michael")
#> [1] "Michaela"      "Michael"      "Michaela"      "Michaelah"      "Michaelyn"
#> [6] "Michael"      "Michaelangelo" "Michaelanthony" "Michaelgabriel"
str_subset(baby2014$name, "michael")
#> [1] "Jamichael"      "Johnmichael"  "Lamichael"      "Jermichael"      "Dmichael"
#> [6] "Jamesmichael"  "Jonmichael"   "Kingmichael"    "Liammichael"     "Carmichael"
#> [11] "Seanmichael"   "Gianmichael"  "Jmichael"
str_subset(baby2014$name, "(M|m)ichael") # contains Michael or michael
#> [1] "Michaela"      "Michael"      "Michaela"      "Michaelah"      "Michaelyn"
#> [6] "Michael"      "Jamichael"    "Michaelangelo" "Johnmichael"    "Lamichael"
#> [11] "Jermichael"    "Michaelanthony" "Dmichael"      "Jamesmichael"   "Jonmichael"
#> [16] "Kingmichael"  "Liammichael"  "Carmichael"    "Michaelgabriel" "Seanmichael"
#> [21] "Gianmichael"  "Jmichael"
str_subset(baby2014$name, "(M|m)ichael$") # ends in Michael or michael
#> [1] "Michael"      "Michael"      "Jamichael"    "Johnmichael"    "Lamichael"
#> [6] "Jermichael"    "Dmichael"      "Jamesmichael" "Jonmichael"     "Kingmichael"
#> [11] "Liammichael"  "Carmichael"    "Seanmichael"  "Gianmichael"    "Jmichael"
str_subset(baby2014$name, "^ (M|m)ichael") # begins with Michael or michael
#> [1] "Michaela"      "Michael"      "Michaela"      "Michaelah"      "Michaelyn"
#> [6] "Michael"      "Michaelangelo" "Michaelanthony" "Michaelgabriel"
str_subset(baby2014$name, "^ (M|m)ichael$") # begins and ends in Michael or michael
#> [1] "Michael" "Michael"

#- find matching records in database
baby2014 %>% filter(str_detect(name, pattern="^ (M|m)ichael$"))
#> # A tibble: 2 × 5
#>   year sex   name      n      prop
#>   <dbl> <chr> <chr> <int>    <dbl>
#> 1  2014   F Michael    35 1.805e-05
#> 2  2014   M Michael 15323 7.549e-03
```

## 6 Examples

### 6.1 Baseball 3,000 Hit Club

```
library(rvest)
url = 'http://en.wikipedia.org/wiki/3,000_hit_club'
hits = read_html(url) %>% html_node("table.wikitable.sortable") %>% html_table()
# head(hits)
```

#### 6.1.1 Fix the player's name (Player)

```
#- Uppercase letter followed by multiple lowercase (last name)
# then comma and space
# then uppercase letter followed by multiple lowercase (first name)
str_extract(hits$Player, "[A-Z][a-z]+, [A-Z][a-z]+")
#> [1] "Rose, Pete"      "Cobb, Ty"        "Aaron, Hank"     "Musial, Stan"
#> [5] "Speaker, Tris"   "Jeter, Derek"    "Wagner, Honus"   "Yastrzemski, Carl"
#> [9] "Molitor, Paul"   "Collins, Eddie"  "Mays, Willie"    "Murray, Eddie"
#> [13] "Lajoie, Nap"     NA                "Brett, George"   "Waner, Paul"
#> [17] "Yount, Robin"    "Gwynn, Tony"     "Rodriguez, Alex" "Winfield, Dave"
#> [21] "Biggio, Craig"   "Henderson, Rickey" "Carew, Rod"      "Suzuki, Ichiro"
#> [25] "Brock, Lou"      "Palmeiro, Rafael" "Anson, Cap"      "Boggs, Wade"
#> [29] "Kaline, Al"      "Clemente, Roberto"
```

Close, but that missed Cal Ripken Jr. So here is another option

```
#- Uppercase letter followed by *anything* (using .)
# then comma and space
# then uppercase letter followed by multiple lowercase (first name)
str_extract(hits$Player, "[A-Z].+, [A-Z][a-z]+")
#> [1] "Rose, Pete"      "Cobb, Ty"        "Aaron, Hank"     "Musial, Stan"
#> [5] "Speaker, Tris"   "Jeter, Derek"    "Wagner, Honus"   "Yastrzemski, Carl"
#> [9] "Molitor, Paul"   "Collins, Eddie"  "Mays, Willie"    "Murray, Eddie"
#> [13] "Lajoie, Nap"     "Ripken Jr., Cal" "Brett, George"   "Waner, Paul"
#> [17] "Yount, Robin"    "Gwynn, Tony"     "Rodriguez, Alex" "Winfield, Dave"
#> [21] "Biggio, Craig"   "Henderson, Rickey" "Carew, Rod"      "Suzuki, Ichiro"
#> [25] "Brock, Lou"      "Palmeiro, Rafael" "Anson, Cap"      "Boggs, Wade"
#> [29] "Kaline, Al"      "Clemente, Roberto"
```

#### 6.1.2 Remove non-numeric from Hits

```
#- We used the readr function parse_number()
parse_number(hits$Hits)
#> [1] 4256 4191 3771 3630 3514 3465 3430 3419 3319 3314 3283 3255 3252 3184 3154 3152 31
#> [18] 3141 3115 3110 3060 3055 3053 3030 3023 3020 3011 3010 3007 3000

#- but we could use a combination of str_extract and str_replace
str_extract(hits$Hits, pattern='\\d+,*\\d+') %>%
  str_replace_all(pattern=',', replacement='') %>%
```



```
as.numeric()
#> [1] 4256 4191 3771 3630 3514 3465 3430 3419 3319 3314 3283 3255 3252 3184 3154 3152 31
#> [18] 3141 3115 3110 3060 3055 3053 3030 3023 3020 3011 3010 3007 3000
```

### 6.1.3 Extract the Date of the 3,000th hit

```
#- We used the pattern in the number of characters
str_sub(hits$Date, 9, 18)
#> [1] "1978-05-05" "1921-08-19" "1970-05-17" "1958-05-13" "1925-05-17" "2011-07-09"
#> [7] "1914-06-09" "1979-09-12" "1996-09-16" "1925-06-03" "1970-07-18" "1995-06-30"
#> [13] "1914-09-27" "2000-04-15" "1992-09-30" "1942-06-19" "1992-09-09" "1999-08-06"
#> [19] "2015-06-19" "1993-09-16" "2007-06-28" "2001-10-07" "1985-08-04" "2016-08-07"
#> [25] "1979-08-13" "2005-07-15" "1897-07-18" "1999-08-07" "1974-09-24" "1972-09-30"

#- But we could just extract with regex
str_extract(hits$Date, pattern='[A-Z][a-z]+ \\d{1,2}, \\d{4}')
```

#> [1]	"May 5, 1978"	"August 19, 1921"	"May 17, 1970"	"May 13, 1958"
#> [5]	"May 17, 1925"	"July 9, 2011"	"June 9, 1914"	"September 12, 1979"
#> [9]	"September 16, 1996"	"June 3, 1925"	"July 18, 1970"	"June 30, 1995"
#> [13]	"September 27, 1914"	"April 15, 2000"	"September 30, 1992"	"June 19, 1942"
#> [17]	"September 9, 1992"	"August 6, 1999"	"June 19, 2015"	"September 16, 1993"
#> [21]	"June 28, 2007"	"October 7, 2001"	"August 4, 1985"	"August 7, 2016"
#> [25]	"August 13, 1979"	"July 15, 2005"	"July 18, 1897"	"August 7, 1999"
#> [29]	"September 24, 1974"	"September 30, 1972"		

```
str_extract(hits$Date, pattern='(18|19|20)\\d{2}-\\d{2}-\\d{2}')
```

#> [1]	"1978-05-05"	"1921-08-19"	"1970-05-17"	"1958-05-13"	"1925-05-17"	"2011-07-09"
#> [7]	"1914-06-09"	"1979-09-12"	"1996-09-16"	"1925-06-03"	"1970-07-18"	"1995-06-30"
#> [13]	"1914-09-27"	"2000-04-15"	"1992-09-30"	"1942-06-19"	"1992-09-09"	"1999-08-06"
#> [19]	"2015-06-19"	"1993-09-16"	"2007-06-28"	"2001-10-07"	"1985-08-04"	"2016-08-07"
#> [25]	"1979-08-13"	"2005-07-15"	"1897-07-18"	"1999-08-07"	"1974-09-24"	"1972-09-30"

### 6.1.4 Get range of seasons played (Seasons)

```
#- We used the str_replace() to remove strange symbols
str_replace_all(hits$Seasons, str_sub(hits$Seasons[1], 5, 7), '-')
#> [1] "1963-1986" "1905-1928" "1954-1976"
#> [4] "1941-1944, 1946-1963" "1907-1928" "1995-2014"
#> [7] "1897-1917" "1961-1983" "1978-1998"
#> [10] "1906-1930" "1951-1952, 1954-1973" "1977-1997"
#> [13] "1896-1916" "1981-2001" "1973-1993"
#> [16] "1926-1945" "1974-1993" "1982-2001"
#> [19] "1994-2013, 2015-2016" "1973-1995" "1988-2007"
#> [22] "1979-2003" "1967-1985" "2001-present"
#> [25] "1961-1979" "1986-2005" "1871-1897"
#> [28] "1982-1999" "1953-1974" "1955-1972"

#- but we can also use str_extract_all(). From this it
# not be too difficult to get the
# *number of seasons* played
```

```

years = '(18|19|20)\\d{2}'
str_extract_all(hits$Seasons, pattern=years, simplify=TRUE)
#>      [,1]      [,2]      [,3]      [,4]
#> [1,] "1963" "1986" ""      ""
#> [2,] "1905" "1928" ""      ""
#> [3,] "1954" "1976" ""      ""
#> [4,] "1941" "1944" "1946" "1963"
#> [5,] "1907" "1928" ""      ""
#> [6,] "1995" "2014" ""      ""
#> [7,] "1897" "1917" ""      ""
#> [8,] "1961" "1983" ""      ""
#> [9,] "1978" "1998" ""      ""
#> [10,] "1906" "1930" ""      ""
#> [11,] "1951" "1952" "1954" "1973"
#> [12,] "1977" "1997" ""      ""
#> [13,] "1896" "1916" ""      ""
#> [14,] "1981" "2001" ""      ""
#> [15,] "1973" "1993" ""      ""
#> [16,] "1926" "1945" ""      ""
#> [17,] "1974" "1993" ""      ""
#> [18,] "1982" "2001" ""      ""
#> [19,] "1994" "2013" "2015" "2016"
#> [20,] "1973" "1995" ""      ""
#> [21,] "1988" "2007" ""      ""
#> [22,] "1979" "2003" ""      ""
#> [23,] "1967" "1985" ""      ""
#> [24,] "2001" ""      ""      ""
#> [25,] "1961" "1979" ""      ""
#> [26,] "1986" "2005" ""      ""
#> [27,] "1871" "1897" ""      ""
#> [28,] "1982" "1999" ""      ""
#> [29,] "1953" "1974" ""      ""
#> [30,] "1955" "1972" ""      ""

```

### 6.1.5 What type of hit was the 3,000th

```

str_extract(hits$`3,000th hit`, "Single|Double|Triple|Home run")
#> [1] "Single" "Single" "Single" "Double" "Single" "Home run" "Double"
#> [8] "Single" "Triple" "Single" "Single" "Single" "Double" "Single"
#> [15] "Single" "Single" "Single" "Single" "Home run" "Single" "Single"
#> [22] "Double" "Single" "Triple" "Single" "Double" "Single" "Home run"
#> [29] "Double" "Double"

```