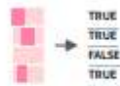# Work with strings with stringr :: **CHEAT SHEET**

stringr

The **stringr** package provides a set of internally consistent tools for working with character strings, i.e. sequences of characters surrounded by quotation marks.
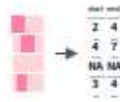
## Detect Matches

**str_detect**(string, **pattern**) Detect the presence of a pattern match in a string. *str_detect(fruit, "a")*

**str_which**(string, **pattern**) Find the indexes of strings that contain a pattern match. *str_which(fruit, "a")*

**str_count**(string, **pattern**) Count the number of matches in a string. *str_count(fruit, "a")*

**str_locate**(string, **pattern**) Locate the positions of pattern matches in a string. Also **str_locate_all**. *str_locate(fruit, "a")*
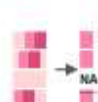
## Subset Strings

**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**) Return only the strings that contain a pattern match. *str_subset(fruit, "b")*

**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)*
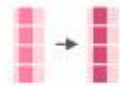
**str_trunc**(string, width, side = c("right", "left", "center"), ellipsis = "...") Truncate the width of strings, replacing content with ellipsis. *str_trunc(fruit, 3)*

**str_trim**(string, side = c("both", "left", "right")) Trim whitespace from the start and/or end of a string. *str_trim(fruit)*
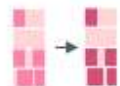
## Mutate Strings

**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. *str_replace(fruit, "a", "-")*
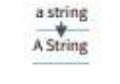
**str_replace_all**(string, **pattern**, replacement) Replace all matched patterns in each string. *str_replace_all(fruit, "a", "-")*

**str_to_lower**(string, locale = "en")[1] Convert strings to lower case. *str_to_lower(sentences)*

**str_to_upper**(string, locale = "en")[1] Convert strings to upper case. *str_to_upper(sentences)*

**str_to_title**(string, locale = "en")[1] Convert strings to title case. *str_to_title(sentences)*

## Join and Split

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

**str_c**(..., sep = "", **collapse = NULL**) Collapse a vector of strings into a single string. *str_c(letters, collapse = "")*

**str_dup**(string, times) Repeat strings times times. *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. *str_split_fixed(fruit, " ", n=2)*

**glue::glue**(..., .sep = "", .envir = parent.frame(), .open = "{", .close = "}") Create a string from strings and {expressions} to evaluate. *glue::glue("Pi is {pi}")*

**glue::glue_data**(.x, ..., .sep = "", .envir = parent.frame(), .open = "{", .close = "}") Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. *glue::glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")*

## 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. *x[str_order(x)]*

**str_sort**(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)[1] Sort a character vector. *str_sort(x)*

## Helpers

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

**str_view**(string, **pattern**, match = NA) View HTML rendering of first regex match in each string. *str_view(fruit, "[aeiou]")*

**str_view_all**(string, **pattern**, match = NA) View HTML rendering of all regex matches. *str_view_all(fruit, "[aeiou]")*

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

# Need to Know

## 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") | abc ABC 123 .!?\(){} |
| \\. | \. | . | see("\\.") | abc ABC 123 .!?\(){} |
| \\! | \! | ! | see("\\!") | abc ABC 123 .!?\(){} |
| \\? | \? | ? | see("\\?") | abc ABC 123 .!?\(){} |
| \\\\ | \\ | \ | see("\\\\") | abc ABC 123 .!?\(){} |
| \\( | \( | ( | see("\\(") | abc ABC 123 .!?\(){} |
| \\) | \) | ) | see("\\)") | abc ABC 123 .!?\(){} |
| \\{ | \{ | { | see("\\{") | abc ABC 123 .!?\(){} |
| \\} | \} | } | see("\\}") | abc ABC 123 .!?\(){} |
| \\n | \n | new line (return) | see("\\n") | abc ABC 123 .!?\(){} |
| \\t | \t | tab | see("\\t") | abc ABC 123 .!?\(){} |
| \\s | \s | any whitespace (\S for non-whitespaces) | see("\\s") | abc ABC 123 .!?\(){} |
| \\d | \d | any digit (\D for non-digits) | see("\\d") | abc ABC 123 .!?\(){} |
| \\w | \w | any word character (\W for non-word chars) | see("\\w") | abc ABC 123 .!?\(){} |
| \\b | \b | word boundaries | see("\\b") | abc ABC 123 .!?\(){} |
| | [:digit:] | digits | see("[:digit:]") | abc ABC 123 .!?\(){} |
| | [:alpha:] | letters | see("[:alpha:]") | abc ABC 123 .!?\(){} |
| | [:lower:] | lowercase letters | see("[:lower:]") | abc ABC 123 .!?\(){} |
| | [:upper:] | uppercase letters | see("[:upper:]") | abc ABC 123 .!?\(){} |
| | [:alnum:] | letters and numbers | see("[:alnum:]") | abc ABC 123 .!?\(){} |
| | [:punct:] | punctuation | see("[:punct:]") | abc ABC 123 .!?\(){} |
| | [:graph:] | letters, numbers, and punctuation | see("[:graph:]") | abc ABC 123 .!?\(){} |
| | [:space:] | space characters (i.e. \s) | see("[:space:]") | abc ABC 123 .!?\(){} |
| | [:blank:] | space and tab (but not new line) | see("[:blank:]") | abc ABC 123 .!?\(){} |
| | . | every character except a new line | see(".") | abc ABC 123 .!?\(){} |

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

[: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 |
| y z | Y 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 | |
|---|---|---|---|
| ^ | start of string | anchor("^a") | aaa |
| $ | end of string | anchor("a$") | aaa |

## LOOK AROUNDS

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

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

## QUANTIFIERS

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

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

# Base R
## Cheat Sheet

## Getting Help

### Accessing the help files

**?mean**
Get help of a particular function.
**help.search('weighted mean')**
Search the help files for a word or phrase.
**help(package = 'dplyr')**
Find help for a package.

### More about an object

**str(iris)**
Get a summary of an object's structure.
**class(iris)**
Find the class an object belongs to.

## Using Packages

**install.packages('dplyr')**
Download and install a package from CRAN.

**library(dplyr)**
Load the package into the session, making all its functions available to use.

**dplyr::select**
Use a particular function from a package.

**data(iris)**
Load a built-in dataset into the environment.

## Working Directory

**getwd()**
Find the current working directory (where inputs are found and outputs are sent).

**setwd('C://file/path')**
Change the current working directory.

**Use projects in RStudio to set the working directory to the folder you are working in.**

## Vectors

### Creating Vectors

| | | |
|---|---|---|
| c(2, 4, 6) | 2 4 6 | Join elements into a vector |
| 2:6 | 2 3 4 5 6 | An integer sequence |
| seq(2, 3, by=0.5) | 2.0 2.5 3.0 | A complex sequence |
| rep(1:2, times=3) | 1 2 1 2 1 2 | Repeat a vector |
| rep(1:2, each=3) | 1 1 1 2 2 2 | Repeat elements of a vector |

### Vector Functions

**sort(x)**
Return x sorted.
**table(x)**
See counts of values.

**rev(x)**
Return x reversed.
**unique(x)**
See unique values.

### Selecting Vector Elements

#### By Position

| | |
|---|---|
| x[4] | The fourth element. |
| x[-4] | All but the fourth. |
| x[2:4] | Elements two to four. |
| x[-(2:4)] | All elements except two to four. |
| x[c(1, 5)] | Elements one and five. |

#### By Value

| | |
|---|---|
| x[x == 10] | Elements which are equal to 10. |
| x[x < 0] | All elements less than zero. |
| x[x %in% c(1, 2, 5)] | Elements in the set 1, 2, 5. |

#### Named Vectors

| | |
|---|---|
| x['apple'] | Element with name 'apple'. |

## Programming

### For Loop

```
for (variable in sequence){
    Do something
}
```

#### Example

```
for (i in 1:4){
    j <- i + 10
    print(j)
}
```

### While Loop

```
while (condition){
    Do something
}
```

#### Example

```
while (i < 5){
    print(i)
    i <- i + 1
}
```

### If Statements

```
if (condition){
    Do something
} else {
    Do something different
}
```

#### Example

```
if (i > 3){
    print('Yes')
} else {
    print('No')
}
```

### Functions

```
function_name <- function(var){
    Do something
    return(new_variable)
}
```

#### Example

```
square <- function(x){
    squared <- x*x
    return(squared)
}
```

### Reading and Writing Data

Also see the **readr** package.

| Input | Ouput | Description |
|---|---|---|
| df <- read.table('file.txt') | write.table(df, 'file.txt') | Read and write a delimited text file. |
| df <- read.csv('file.csv') | write.csv(df, 'file.csv') | Read and write a comma separated value file. This is a special case of read.table/write.table. |
| load('file.RData') | save(df, file = 'file.Rdata') | Read and write an R data file, a file type special for R. |

| Conditions | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | a == b | Are equal | a > b | Greater than | a >= b | Greater than or equal to | is.na(a) | Is missing |
| | a != b | Not equal | a < b | Less than | a <= b | Less than or equal to | is.null(a) | Is null |

# Types

Converting between common data types in R. Can always go from a higher value in the table to a lower value.

| | | |
|---|---|---|
| `as.logical` | TRUE, FALSE, TRUE | Boolean values (TRUE or FALSE). |
| `as.numeric` | 1, 0, 1 | Integers or floating point numbers. |
| `as.character` | '1', '0', '1' | Character strings. Generally preferred to factors. |
| `as.factor` | '1', '0', '1', levels: '1', '0' | Character strings with preset levels. Needed for some statistical models. |

# Maths Functions

| | | | |
|---|---|---|---|
| `log(x)` | Natural log. | `sum(x)` | Sum. |
| `exp(x)` | Exponential. | `mean(x)` | Mean. |
| `max(x)` | Largest element. | `median(x)` | Median. |
| `min(x)` | Smallest element. | `quantile(x)` | Percentage quantiles. |
| `round(x, n)` | Round to n decimal places. | `rank(x)` | Rank of elements. |
| `signif(x, n)` | Round to n significant figures. | `var(x)` | The variance. |
| `cor(x, y)` | Correlation. | `sd(x)` | The standard deviation. |

# Variable Assignment

```
> a <- 'apple'
> a
[1] 'apple'
```
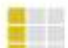
# The Environment

| | |
|---|---|
| `ls()` | List all variables in the environment. |
| `rm(x)` | Remove x from the environment. |
| `rm(list = ls())` | Remove all variables from the environment. |

**You can use the environment panel in RStudio to browse variables in your environment.**

# Matrices

`m <- matrix(x, nrow = 3, ncol = 3)`
Create a matrix from x.

`m[2, ]` - Select a row

`m[ , 1]` - Select a column

`m[2, 3]` - Select an element

`t(m)`
Transpose
`m %*% n`
Matrix Multiplication
`solve(m, n)`
Find x in: m * x = n

# Lists

`l <- list(x = 1:5, y = c('a', 'b'))`
A list is a collection of elements which can be of different types.

| `l[[2]]` | `l[1]` | `l$x` | `l['y']` |
|---|---|---|---|
| Second element of l. | New list with only the first element. | Element named x. | New list with only element named y. |

# Data Frames

`df <- data.frame(x = 1:3, y = c('a', 'b', 'c'))`
A special case of a list where all elements are the same length.

| x | y |
|---|---|
| 1 | a |
| 2 | b |
| 3 | c |

**List subsetting**

`df$x`    `df[[2]]`

*Understanding a data frame*

| `View(df)` | See the full data frame. |
|---|---|
| `head(df)` | See the first 6 rows. |

**Matrix subsetting**

`df[ , 2]`

`df[2, ]`

`df[2, 2]`

| `nrow(df)` | Number of rows. |
|---|---|
| `ncol(df)` | Number of columns. |
| `dim(df)` | Number of columns and rows. |

`cbind` - Bind columns.

`rbind` - Bind rows.

# Strings

| | |
|---|---|
| `paste(x, y, sep = ' ')` | Join multiple vectors together. |
| `paste(x, collapse = ' ')` | Join elements of a vector together. |
| `grep(pattern, x)` | Find regular expression matches in x. |
| `gsub(pattern, replace, x)` | Replace matches in x with a string. |
| `toupper(x)` | Convert to uppercase. |
| `tolower(x)` | Convert to lowercase. |
| `nchar(x)` | Number of characters in a string. |

# Factors

`factor(x)`
Turn a vector into a factor. Can set the levels of the factor and the order.

`cut(x, breaks = 4)`
Turn a numeric vector into a factor by 'cutting' into sections.

# Statistics

`lm(y ~ x, data=df)`
Linear model.

`glm(y ~ x, data=df)`
Generalised linear model.

`summary`
Get more detailed information out a model.

`t.test(x, y)`
Perform a t-test for difference between means.

`pairwise.t.test`
Perform a t-test for paired data.

`prop.test`
Test for a difference between proportions.

`aov`
Analysis of variance.

# Distributions

| | Random Variates | Density Function | Cumulative Distribution | Quantile |
|---|---|---|---|---|
| Normal | `rnorm` | `dnorm` | `pnorm` | `qnorm` |
| Poisson | `rpois` | `dpois` | `ppois` | `qpois` |
| Binomial | `rbinom` | `dbinom` | `pbinom` | `qbinom` |
| Uniform | `runif` | `dunif` | `punif` | `qunif` |

# Plotting

`plot(x)`
Values of x in order.

`plot(x, y)`
Values of x against y.

`hist(x)`
Histogram of x.

# Dates

# par() Graphical Parameters

Visual cheat sheet for some plot parameters in R. See ?par for more information.

## Symbol Styles

### pch | Point Types

| | | | |
|---|---|---|---|
| ○ 1 | ◩ 14 | | |
| △ 2 | ■ 15 | | |
| + 3 | ● 16 | | |
| × 4 | ▲ 17 | | |
| ◇ 5 | ◆ 18 | | |
| ▽ 6 | ● 19 | | |
| ⊠ 7 | ● 20 | | |
| ✳ 8 | ◉ 21 | | |
| ⟠ 9 | ▢ 22 | | |
| ⊕ 10 | ◈ 23 | | |
| ⧖ 11 | ▲ 24 | | |
| ⊞ 12 | ▽ 25 | | |
| ⊗ 13 | you can also use any character | | |

### lty | Line Types

1, 2, 3, 4, 5, 6

### lwd | Line Width

.1, .25, .5, 1, 3, 6

## Figures Arrangement

### mfrow | Multiple Figures by Row

**2,3**

| 1 | 2 | 3 |
|---|---|---|
| 4 | 5 | 6 |

Also available **mfcol** for multiple figures by column

## Axes

### lab | Tick Placement

**10,10** **1,10**

**10,1** **2,2**

### tck | Tick Length

**-0.1** **0.1** **1**

### bty | Box Type

**'o'** **'l'** **'7'**

**'c'** **'u'** **']'**

## Text and Labels

### family, font | Typeface and Font Style

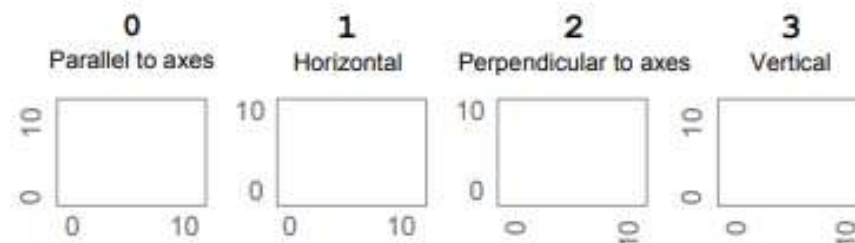| family: mono font: 1 | family: serif font: 1 | family: sans font: 1 |
|---|---|---|
| **family: mono font: 2** | **family: serif font: 2** | **family: sans font: 2** |
| *family: mono font: 3* | *family: serif font: 3* | *family: sans font: 3* |
| ***family: mono font: 4*** | ***family: serif font: 4*** | ***family: sans font: 4*** |

Also available: **font.main** (main title), **font.lab** (axis labels), **font.sub** (subtitle)

### las | Label Orientation

| **0** | **1** | **2** | **3** |
|---|---|---|---|
| Parallel to axes | Horizontal | Perpendicular to axes | Vertical |

### ann | Plot Annotation

**TRUE** **FALSE**

Some Title

y-values

x-values

### srt | String Rotation

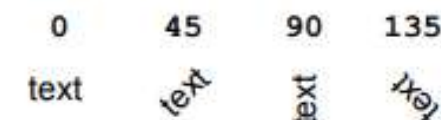| 0 | 45 | 90 | 135 |
|---|---|---|---|
| text | text | text | text |

### lheight | Line Height

**1**

The quick brown fox jumps over the lazy dog and runs away with all the food

**1.5**

The quick brown fox jumps over the lazy dog and runs away with all the food

*Based on Flowing Data's cheat sheet*

# How Big is Your Graph?
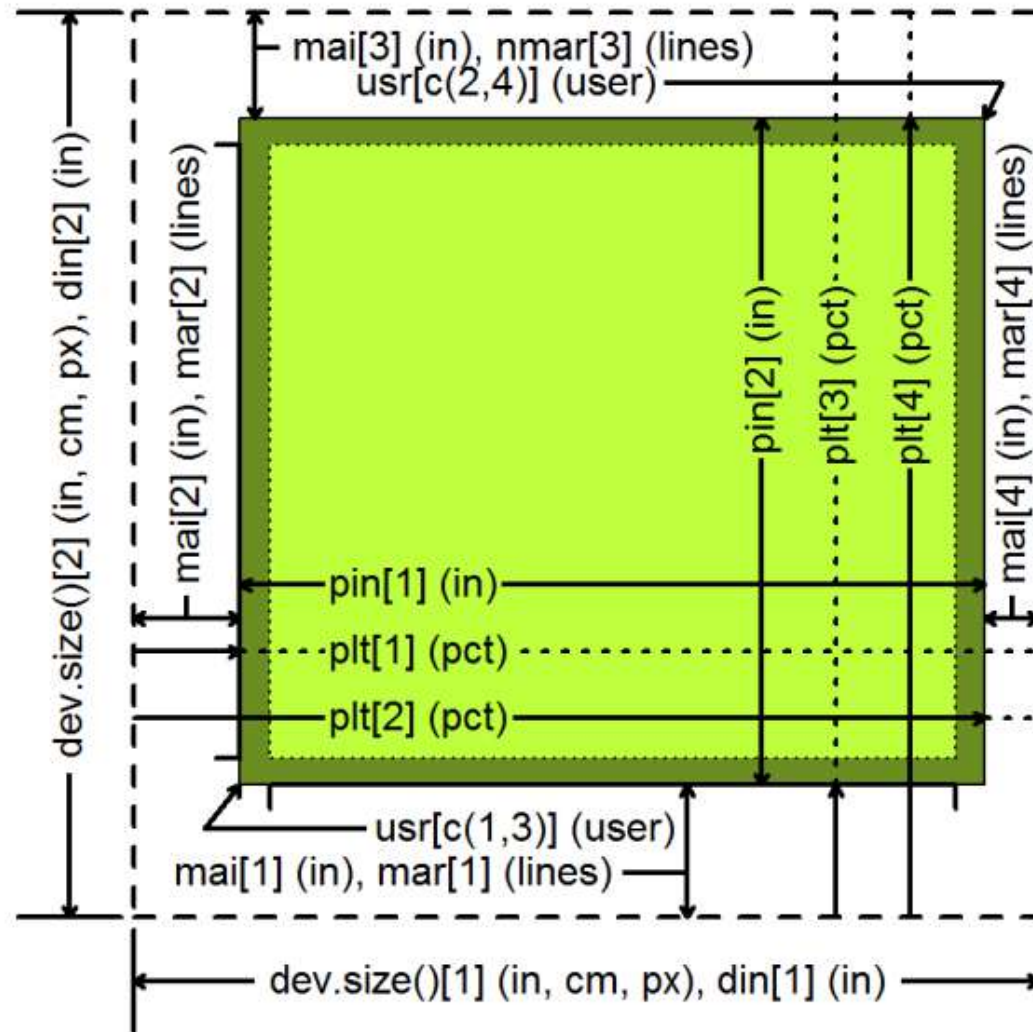
## An R Cheat Sheet

## Introduction

All functions that open a device for graphics will have **height** and **width** arguments to control the size of the graph and a **pointsize** argument to control the relative font size. In **knitr**, you control the size of the graph with the chunk options, **fig.width** and **fig.height**. This sheet will help you with calculating the size of the graph and various parts of the graph within R.



## Your graphics device

**dev.size()** (width, height)
**par("din")** *(r.o.)* (width, height) in inches

Both the **dev.size** function and the **din** argument of **par** will tell you the size of the graphics device. The **dev.size** function will report the size in

1. inches (**units="in"**), the default
2. centimeters (**units="cm"**)
3. pixels (**units="px"**)

Like several other **par** arguments, **din** is read only *(r.o.)* meaning that you can ask its current value (**par("din")**) but you cannot change it (**par(din=c(5,7))** will fail).

## Your plot margins

**par("mai")** (bottom, left, top, right) in inches
**par("mar")** (bottom, left, top, right) in lines

Margins provide you space for your axes, axis, labels, and titles.

A "line" is the amount of vertical space needed for a line of text.

If your graph has no axes or titles, you can remove the margins (and maximize the plotting region) with

**par(mar=rep(0,4))**

## Your plotting region

**par("pin")** (width, height) in inches
**par("plt")** (left, right, bottom, top) in pct

The **pin** argument **par** gives you the size of the plotting region (the size of the device minus the size of the margins) in inches.

The **plt** argument gives you the percentage of the device from the left/bottom edge up to the left edge of the plotting region, the right edge, the bottom edge, and the top edge. The first and third values are equivalent to the percentage of space devoted to the left and bottom margins. Subtract the second and fourth values from 1 to get the percentage of space devoted to the right and top margins.

## Your x-y coordinates

**par("usr")** (xmin, ymin, xmax, ymax)

Your x-y coordinates are the values you use when plotting your data. This normally is not the same as the values you specified with the **xlim** and **ylim** arguments in **plot**. By default, R adds an extra 4% to the plotting range (see the dark green region on the figure) so that points right up on the edges of your plot do not get partially clipped. You can override this by setting **xaxs="i"** and/or the **yaxs="i"** in **par**.

Run **par("usr")** to find the minimum X value, the maximum X value, the minimum Y value, and the maximum Y value. If you assign new values to **usr**, you will update the x-y coordinates to the new values.

## Getting a square graph

**par("pty")**

You can produce a square graph manually by setting the width and height to the same value and setting the margins so that the sum of the top and bottom margins equal the sum of the left and right margins. But a much easier way is to specify **pty="s"**, which adjusts the margins so that the size of the plotting region is always square, even if you resize the graphics window.

## Converting units

For many applications, you need to be able to translate user coordinates to pixels or inches. There are some cryptic shortcuts, but the simplest way is to get the range in user coordinates and measure the proportion of the graphics device devoted to the plotting region.

**user.range <- par("usr")[c(2,4)] - par("usr")[c(1,3)]**

**region.pct <- par("plt")[c(2,4)] - par("plt")[c(1,3)]**

**region.px <- dev.size(units="px") * region.pct**

**px.per.xy <- region.px / user.range**

To convert a horizontal or distance from the x-coordinate value to pixels, multiply by **px.per.xy[1]**. To convert a vertical distance, multiply by **region.px.per.xy[2]**. To convert a diagonal distance, you need to invoke Pyhthagoras.

**a.px <- x.dist*px.per.xy[1]**
**b.px <- y.dist*px.per.xy[2]**
**c.px <- sqrt(a.px^2+b.px^2)**

To rotate a string to match the slope of a line segment, you need to convert the distances to pixels, calculate the arctangent, and convert from radians to degrees.

**segments(x0, y0, x1, y1)**
**delta.x <- (x1 – x0) * px.per.xy[1]**
**delta.y <- (y1 – y0) * px.per.xy[y]**
**angle.radians <- atan2(delta.y, delta.x)**
**angle.degrees <- angle.radians * 180 / pi**
**text(x1, y1, "TEXT", srt=angle.degrees)**

## Panels

**par("fig")** (width, height) in pct
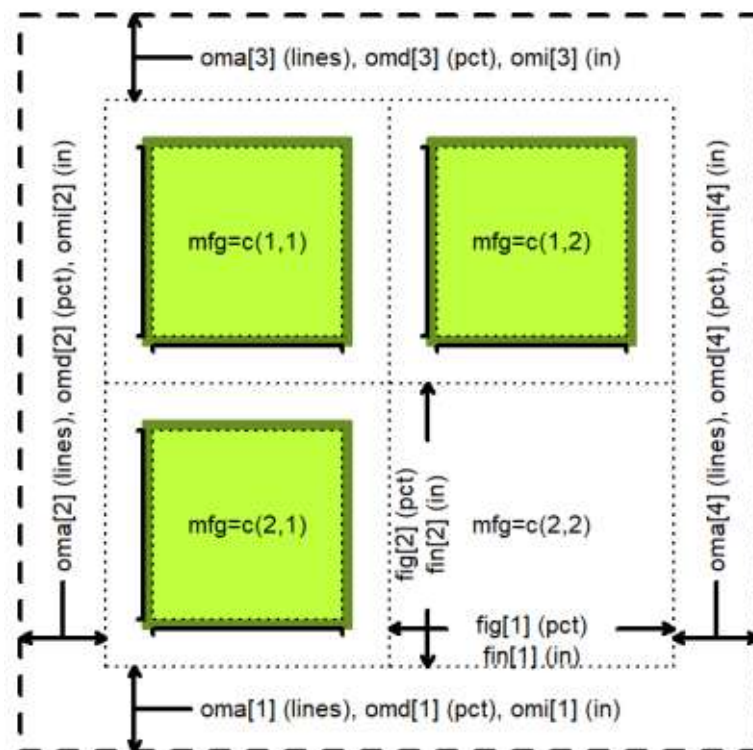**par("fin")** (width, height) in inches

If you display multiple plots within a single graphics window (e.g., with the **mfrow** or **mfcol** arguments of **par** or with the **layout** function), then the **fig** and **fin** arguments will tell you the size of the current subplot window in percent or inches, respectively.

**par("oma")** (bottom, left, top, right) in lines
**par("omd")** (bottom, left, top, right) in pct
**par("omi")** (bottom, left, top, right) in inches

Each subplot will have margins specified by **mai** or **mar**, but no outer margin around the entire set of plots, unless you specify them using **oma**, **omd**, or **omi**. You can place text in the outer margins using the **mtext** function with the argument **outer=TRUE**.

**par("mfg")** (r, c) or (r, c, maxr, maxc)

The **mfg** argument of **par** will allow you to jump to a subplot in a particular row and column. If you query with **par("mfg")**, you will get the current row and column followed by the maximum row and column.

## Character and string sizes

### strheight()

The **strheight** functions will tell you the height of a specified string in inches (**units="inches"**), x-y user coordinates (**units="user"**) or as a percentage of the graphics device (**units="figure"**).

For a single line of text, **strheight** will give you the height of the letter "M". If you have a string with one of more linebreaks ("\n"), the **strheight** function will measure the height of the letter "M" plus the height of one or more additional lines. The height of a line is dependent on the line spacing, set by the **lheight** argument of **par**. The default line height (**lheight=1**), corresponding to single spaced lines, produces a line height roughly 1.5 times the height of "M".

### strwidth()

The **strwidth** function will produce different widths to individual characters, representing the proportional spacing used by most fonts (a "W" using much more space than an "i"). For the width of a string, the **strwidth** function will sum up the lengths of the individual characters in the string.
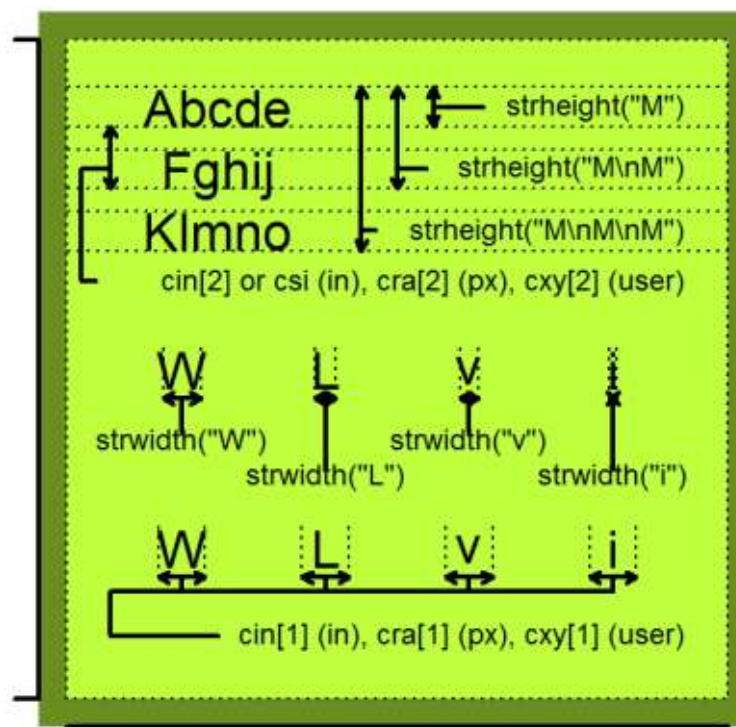
**par("cin")** (r.o.) (width, height) in inches
**par("csi")** (r.o.) height in inches
**par("cra")** (r.o.) (width, height) in pixels
**par("cxy")** (r.o.) (width, height) in xy coordinates

The single value returned by the **csi** argument of **par** gives you the height of a line of text in inches. The second of the two values returned by **cin**, **cra**, and **cxy** gives you the height of a line, in inches, pixels, or xy (user) coordinates.

The first of the two values returned by the **cin**, **cra**, and **cxy** arguments to **par** gives you the approximate width of a single character, in inches, pixels, or xy (user) coordinates. The width, very slightly smaller than the actual width of the letter "W", is a rough estimate at best and ignores the variable with of individual letters.

These values are useful, however, in providing fast ratios of the relative sizes of the differing units of measure

**px.per.in <- par("cra") / par("cin")**
**px.per.xy <- par("cra") / par("cxy")**
**xy.per.in <- par("cxy") / par("cin")**

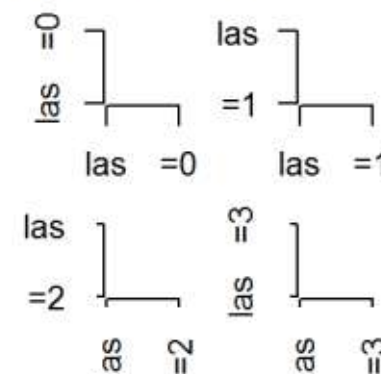## If your fonts are too big or too small

Fixing this takes a bit of trial and error.

1. Specify a larger/smaller value for the **pointsize** argument when you open your graphics device.

2. Trying opening your graphics device with different values for **height** and **width**. Fonts that look too big might be better proportioned in a larger graphics window.

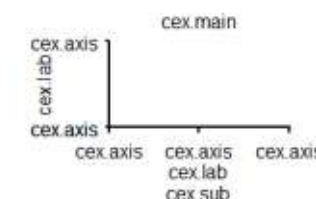3. Use the **cex** argument to increase or decrease the relative size of your fonts.

## If your axes don't fit

There are several possible solutions.

1. You can assign wider margins using the **mar** or **mai** argument in **par**.

2. You can change the orientation of the axis labels with **las**. Choose among
   a. **las=0** both axis labels parallel
   b. **las=1** both axis labels horizontal
   c. **las=2** both axis labels perpendicular
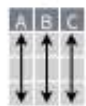   d. **las=3** both axis labels vertical.

3. change the relative size of the font
   a. **cex.axis** for the tick mark labels.
   b. **cex.lab** for **xlab** and **ylab**.
   c. **cex.main** for the main title.
   d. **cex.sub** for the subtitle.

# Data tidying with tidyr : : CHEAT SHEET

**Tidy data** is a way to organize tabular data in a consistent data structure across packages. A table is tidy if:
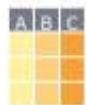
Each **variable** is in its own **column**

&

Each **observation**, or **case**, is in its own row

Access **variables** as **vectors**

Preserve **cases** in vectorized operations

## Tibbles

### AN ENHANCED DATA FRAME

Tibbles are a table format provided by the **tibble** package. They inherit the data frame class, but have improved behaviors:

- **Subset** a new tibble with ], a vector with [[ and $.
- **No partial matching** when subsetting columns.
- **Display** concise views of the data on one screen.

**options**(tibble.print_max = n, tibble.print_min = m, tibble.width = Inf) Control default display settings.

**View()** or **glimpse()** View the entire data set.

### CONSTRUCT A TIBBLE

**tibble(…)** Construct by columns.
tibble(x = 1:3, y = c("a", "b", "c"))

**tribble(…)** Construct by rows.
tribble(~x, ~y,
    1, "a",
    2, "b",
    3, "c")

Both make this tibble

```
A tibble: 3 × 2
      x     y
  <int> <chr>
1     1     a
2     2     b
3     3     c
```

**as_tibble**(x, …) Convert a data frame to a tibble.

**enframe**(x, name = "name", value = "value") Convert a named vector to a tibble. Also **deframe()**.

**is_tibble**(x) Test whether x is a tibble.

## Reshape Data - Pivot data to reorganize values into a new layout.

**pivot_longer**(data, cols, names_to = "name", values_to = "value", values_drop_na = FALSE)

"Lengthen" data by collapsing several columns into two. Column names move to a new names_to column and values to a new values_to column.

pivot_longer(table4a, cols = 2:3, names_to = "year", values_to = "cases")

**pivot_wider**(data, names_from = "name", values_from = "value")

The inverse of pivot_longer(). "Widen" data by expanding two columns into several. One column provides the new column names, the other the values.

pivot_wider(table2, names_from = type, values_from = count)

## Split Cells - Use these functions to split or combine cells into individual, isolated values.

**unite**(data, col, …, sep = "_", remove = TRUE, na.rm = FALSE) Collapse cells across several columns into a single column.

unite(table5, century, year, col = "year", sep = "")

**separate**(data, col, into, sep = "[^[:alnum:]]+", remove = TRUE, convert = FALSE, extra = "warn", fill = "warn", …) Separate each cell in a column into several columns. Also **extract()**.

separate(table3, rate, sep = "/", into = c("cases", "pop"))

**separate_rows**(data, …, sep = "[^[:alnum:].]+", convert = FALSE) Separate each cell in a column into several rows.

separate_rows(table3, rate, sep = "/")

## Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).
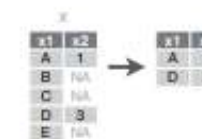
**expand**(data, …) Create a new tibble with all possible combinations of the values of the variables listed in … Drop other variables.
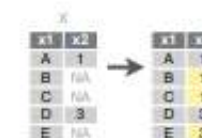expand(mtcars, cyl, gear, carb)

**complete**(data, …, fill = list()) Add missing possible combinations of values of variables listed in … Fill remaining variables with NA.
complete(mtcars, cyl, gear, carb)

## Handle Missing Values

Drop or replace explicit missing values (NA).

**drop_na**(data, …) Drop rows containing NA's in … columns.
drop_na(x, x2)

**fill**(data, …, .direction = "down") Fill in NA's in … columns using the next or previous value.
fill(x, x2)

**replace_na**(data, replace) Specify a value to replace NA in selected columns.
replace_na(x, list(x2 = 2))

# Data Transformation with dplyr :: CHEAT SHEET

**dplyr** functions work with pipes and expect **tidy data**. In tidy data:

Each **variable** is in its own **column**

&

Each **observation**, or **case**, is in its own **row**

**pipes**

x %>% f(y)
becomes f(x, y)

## Summarise Cases

These apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).

**summary function**

**summarise**(.data, ...)
Compute table of summaries.
*summarise(mtcars, avg = mean(mpg))*

**count**(x, ..., wt = NULL, sort = FALSE)
Count number of rows in each group defined by the variables in ... Also **tally**().
*count(iris, Species)*

### VARIATIONS

**summarise_all()** - Apply funs to every column.
**summarise_at()** - Apply funs to specific columns.
**summarise_if()** - Apply funs to all cols of one type.

## Group Cases

Use **group_by()** to create a "grouped" copy of a table. dplyr functions will manipulate each "group" separately and then combine the results.

mtcars %>%

group_by(cyl) %>%

summarise(avg = mean(mpg))

**group_by**(.data, ..., add = FALSE)
Returns copy of table grouped by ...
g_iris <- group_by(iris, Species)

**ungroup**(x, ...)
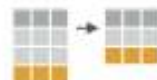Returns ungrouped copy of table.
ungroup(g_iris)

## Manipulate Cases

### EXTRACT CASES

Row functions return a subset of rows as a new table.

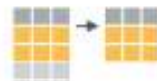**filter**(.data, ...) Extract rows that meet logical criteria. *filter(iris, Sepal.Length > 7)*

**distinct**(.data, ..., .keep_all = FALSE) Remove rows with duplicate values.
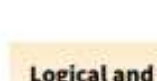*distinct(iris, Species)*

**sample_frac**(tbl, size = 1, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select fraction of rows.
*sample_frac(iris, 0.5, replace = TRUE)*

**sample_n**(tbl, size, replace = FALSE, weight = NULL, .env = parent.frame()) Randomly select size rows. *sample_n(iris, 10, replace = TRUE)*

**slice**(.data, ...) Select rows by position.
*slice(iris, 10:15)*

**top_n**(x, n, wt) Select and order top n entries (by group if grouped data). *top_n(iris, 5, Sepal.Width)*

> **Logical and boolean operators to use with filter()**
>
> | < | <= | is.na() | %in% | \| | xor() |
> |---|---|---|---|---|---|
> | > | >= | !is.na() | ! | & | |
>
> See **?base::logic** and **?Comparison** for help.

### ARRANGE CASES

**arrange**(.data, ...) Order rows by values of a column or columns (low to high), use with **desc()** to order from high to low.
arrange(mtcars, mpg)
arrange(mtcars, desc(mpg))

### ADD CASES

**add_row**(.data, ..., .before = NULL, .after = NULL) Add one or more rows to a table.
*add_row(faithful, eruptions = 1, waiting = 1)*

## Manipulate Variables

### EXTRACT VARIABLES

Column functions return a set of columns as a new vector or table.

**pull**(.data, var = -1) Extract column values as a vector. Choose by name or index.
*pull(iris, Sepal.Length)*

**select**(.data, ...)
Extract columns as a table. Also **select_if()**.
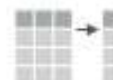*select(iris, Sepal.Length, Species)*

**Use these helpers with select ()**,
*e.g. select(iris, starts_with("Sepal"))*

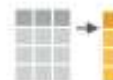| | | |
|---|---|---|
| **contains**(match) | **num_range**(prefix, range) | **:**, e.g. mpg:cyl |
| **ends_with**(match) | **one_of**(...) | **-**, e.g, -Species |
| **matches**(match) | **starts_with**(match) | |

### MAKE NEW VARIABLES

These apply **vectorized functions** to columns. Vectorized funs take vectors as input and return vectors of the same length as output (see back).
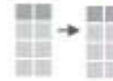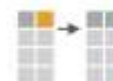
**vectorized function**

**mutate**(.data, ...)
Compute new column(s).
*mutate(mtcars, gpm = 1/mpg)*

**transmute**(.data, ...)
Compute new column(s), drop others.
*transmute(mtcars, gpm = 1/mpg)*

**mutate_all**(.tbl, .funs, ...) Apply funs to every column. Use with **funs()**. Also **mutate_if()**.
*mutate_all(faithful, funs(log(.), log2(.)))*
*mutate_if(iris, is.numeric, funs(log(.)))*

**mutate_at**(.tbl, .cols, .funs, ...) Apply funs to specific columns. Use with **funs()**, **vars()** and the helper functions for select().
*mutate_at(iris, vars( -Species), funs(log(.)))*

**add_column**(.data, ..., .before = NULL, .after = NULL) Add new column(s). Also **add_count()**, **add_tally()**. *add_column(mtcars, new = 1:32)*

**rename**(.data, ...) Rename columns.
*rename(iris, Length = Sepal.Length)*

# Vector Functions

## TO USE WITH MUTATE ()

**mutate()** and **transmute()** apply vectorized functions to columns to create new columns. Vectorized functions take vectors as input and return vectors of the same length as output.

| vectorized function →

### OFFSETS

dplyr::**lag()** - Offset elements by 1
dplyr::**lead()** - Offset elements by -1

### CUMULATIVE AGGREGATES

dplyr::**cumall()** - Cumulative all()
dplyr::**cumany()** - Cumulative any()
    **cummax()** - Cumulative max()
dplyr::**cummean()** - Cumulative mean()
    **cummin()** - Cumulative min()
    **cumprod()** - Cumulative prod()
    **cumsum()** - Cumulative sum()

### RANKINGS

dplyr::**cume_dist()** - Proportion of all values <=
dplyr::**dense_rank()** - rank with ties = min, no gaps
dplyr::**min_rank()** - rank with ties = min
dplyr::**ntile()** - bins into n bins
dplyr::**percent_rank()** - min_rank scaled to [0,1]
dplyr::**row_number()** - rank with ties = "first"

### MATH

**+, -, *, /, ^, %/%, %%** - arithmetic ops
**log(), log2(), log10()** - logs
**<, <=, >, >=, !=, ==** - logical comparisons
dplyr::**between()** - x >= left & x <= right
dplyr::**near()** - safe == for floating point numbers

### MISC

dplyr::**case_when()** - multi-case if_else()
dplyr::**coalesce()** - first non-NA values by element across a set of vectors
dplyr::**if_else()** - element-wise if() + else()
dplyr::**na_if()** - replace specific values with NA
    **pmax()** - element-wise max()
    **pmin()** - element-wise min()
dplyr::**recode()** - Vectorized switch()
dplyr::**recode_factor()** - Vectorized switch() for factors

---

# Summary Functions

## TO USE WITH SUMMARISE ()

**summarise()** applies summary functions to columns to create a new table. Summary functions take vectors as input and return single values as output.

| summary function →

### COUNTS

dplyr::**n()** - number of values/rows
dplyr::**n_distinct()** - # of uniques
    **sum(!is.na())** - # of non-NA's

### LOCATION

    **mean()** - mean, also **mean(!is.na())**
    **median()** - median

### LOGICALS

    **mean()** - Proportion of TRUE's
    **sum()** - # of TRUE's

### POSITION/ORDER

dplyr::**first()** - first value
dplyr::**last()** - last value
dplyr::**nth()** - value in nth location of vector

### RANK

    **quantile()** - nth quantile
    **min()** - minimum value
    **max()** - maximum value

### SPREAD

    **IQR()** - Inter-Quartile Range
    **mad()** - median absolute deviation
    **sd()** - standard deviation
    **var()** - variance

---

# Row Names

Tidy data does not use rownames, which store a variable outside of the columns. To work with the rownames, first move them into a column.

 **rownames_to_column()**
Move row names into col.
a <- rownames_to_column(iris, var = "C")

 **column_to_rownames()**
Move col in row names.
column_to_rownames(a, var = "C")

Also **has_rownames()**, **remove_rownames()**

---

# Combine Tables

## COMBINE VARIABLES



Use **bind_cols()** to paste tables beside each other as they are.

**bind_cols(...)** Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

Use a "**Mutating Join**" to join one table to columns from another, matching values with the rows that they correspond to. Each join retains a different combination of values from the tables.

 **left_join**(x, y, by = NULL, copy=FALSE, suffix=c(".x",".y"),...)
Join matching values from y to x.

 **right_join**(x, y, by = NULL, copy = FALSE, suffix=c(".x",".y"),...)
Join matching values from x to y.

 **inner_join**(x, y, by = NULL, copy = FALSE, suffix=c(".x",".y"),...)
Join data. Retain only rows with matches.

 **full_join**(x, y, by = NULL, copy=FALSE, suffix=c(".x",".y"),...)
Join data. Retain all values, all rows.

 Use **by = c("col1", "col2")** to specify the column(s) to match on.
left_join(x, y, by = "A")

 Use a named vector, **by = c("col1" = "col2")**, to match on columns with different names in each data set.
left_join(x, y, by = c("C" = "D"))

 Use **suffix** to specify suffix to give to duplicate column names.
left_join(x, y, by = c("C" = "D"), suffix = c("1", "2"))

## COMBINE CASES



Use **bind_rows()** to paste tables below each other as they are.

 **bind_rows(..., .id = NULL)**
Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured)

 **intersect(x, y, ...)**
Rows that appear in both x and y.

 **setdiff(x, y, ...)**
Rows that appear in x but not y.

 **union(x, y, ...)**
Rows that appear in x or y. (Duplicates removed). union_all() retains duplicates.

Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

## EXTRACT ROWS



Use a "**Filtering Join**" to filter one table against the rows of another.

 **semi_join(x, y, by = NULL, ...)**
Return rows of x that have a match in y. USEFUL TO SEE WHAT WILL BE JOINED.

 **anti_join(x, y, by = NULL, ...)**
Return rows of x that do not have a match in y. USEFUL TO SEE WHAT WILL NOT BE JOINED.