

Data import With the tidyverse :: CHEAT SHEET

Read Tabular Data with `readr`

```
read_*(file, col_names = TRUE, col_types = NULL, col_select = NULL, id = NULL, locale, n_max = Inf,
skip = 0, na = c("", "NA"), guess_max = min(1000, n_max), show_col_types = TRUE) See ?read_delim
```

A|B|C
1|2|3
4|5|NA

→ A | B | C
1 | 2 | 3
4 | 5 | NA

read_delim("file.txt", delim = "|") Read files with any delimiter. If no delimiter is specified, it will automatically guess.
To make file.txt, run: write_file("A|B|C\n1|2|3\n4|5|NA", file = "file.txt")

A|B|C
1|2|3
4|5|NA

→ A | B | C
1 | 2 | 3
4 | 5 | NA

read_csv("file.csv") Read a comma delimited file with period decimal marks.
write_file("A,B,C\n1,2,3\n4,5,NA", file = "file.csv")

A|B|C
1|2|3
4|5|NA

→ A | B | C
1 | 2 | 3
4 | 5 | NA

A|B|C
1|2|3
4|5|NA

→ A | B | C
1 | 2 | 3
4 | 5 | NA

read_csv2("file2.csv") Read semicolon delimited files with comma decimal marks.
write_file("A;B;C\n1;2;3\n4;5;NA", file = "file2.csv")

A|B|C
1|2|3
4|5|NA

→ A | B | C
1 | 2 | 3
4 | 5 | NA

read_tsv("file.tsv") Read a tab delimited file. Also **read_table()**.
read_fwf("file.tsv", fwf_widths(c(2, 2, NA))) Read a fixed width file.
write_file("A\tB\tC\n1\t2\t3\n4\t5\tNA", file = "file.tsv")

USEFUL READ ARGUMENTS

A | B | C
1 | 2 | 3
4 | 5 | NA

No header
read_csv("file.csv", col_names = FALSE)

x | y | z
A | B | C

Provide header
read_csv("file.csv",
col_names = c("x", "y", "z"))

A | B | C
1 | 2 | 3
4 | 5 | NA

Read values as missing
read_csv("file.csv", na = c("1"))

→ ████
A | B | C
1 | 2 | 3
4 | 5 | NA

Read multiple files into a single table
read_csv(c("f1.csv", "f2.csv", "f3.csv"),
id = "origin_file")

A|B|C
1|2|3|0

Save Data with `readr`

```
write_*(x, file, na = "NA", append, col_names, quote, escape, eol, num_threads, progress)
```

→ A | B | C
1 | 2 | 3
4 | 5 | NA

write_delim(x, file, delim = "") Write files with any delimiter.

write_csv(x, file) Write a comma delimited file.

write_csv2(x, file) Write a semicolon delimited file.

write_tsv(x, file) Write a tab delimited file.



One of the first steps of a project is to import outside data into R. Data is often stored in tabular formats, like csv files or spreadsheets.

The front page of this sheet shows how to import and save text files into R using `readr`.

The back page shows how to import spreadsheet data from Excel files using `readxl` or Google Sheets using `googlesheets4`.

OTHER TYPES OF DATA

Try one of the following packages to import other types of files:

- haven** - SPSS, Stata, and SAS files
- jsonlite** - json
- xml2** - XML
- httr** - Web APIs
- rvest** - HTML (Web Scraping)
- readr::read_lines()** - text data

Column Specification with `readr`

Column specifications define what data type each column of a file will be imported as. By default readr will generate a column spec when a file is read and output a summary.

spec(x) Extract the full column specification for the given imported data frame.

Skip lines
read_csv("file.csv", skip = 1)
age is an integer

Read a subset of lines
read_csv("file.csv", n_max = 1)

Read values as missing
read_csv("file.csv", na = c("1"))

COLUMN TYPES

Each column type has a function and corresponding string abbreviation.

age is an integer
sex is a character

DEFINE COLUMN SPECIFICATION

Set a default type
read_csv(
file,
col_type = list(default = col_double()))

Use column type or string abbreviation

• col_character() - "c"
• col_factor(levels, ordered = FALSE) - "f"
• col_datetime(format = "") - "T"
• col_date(format = "") - "D"
• col_time(format = "") - "t"
• col_skip() - "_" - "
• col_double() - "d"
• col_number() - "n"

Use a single string of abbreviations

col_types: skip, guess, integer, logical, character
read_csv(
file,
col_type = "?|C")

Import Spreadsheets

with readxl

READ EXCEL FILES

A	B	C	D	E
1	x1	x2	x3	x4
2	x	z	8	NA
3	y	7	9	10

x1	x2	x3	x4	x5
x	NA	z	8	NA
x	z	8	NA	9
y	7	NA	9	10

read_excel(path, sheet = NULL, range = NULL)

Read a .xls or .xlsx file based on the file extension.

See front page for more read arguments. Also

read_xls() and **read_xlsx()**.

`read_excel("excel_file.xls")`

READ SHEETS

read_excel(path, sheet = NULL) Specify which sheet to read by position or name.

`read_excel(path, sheet = 1)`

`read_excel(path, sheet = "s1")`

read_excel(path, sheet = "s1") Get a vector of sheet names.

`excel_sheets("excel_file.xls")`

guess_max

To guess a column type, `read_excel()` looks at the first 1000 rows of data. Increase with the

guess_max argument.

`read_excel(path, guess_max = Inf)`

Set all columns to same type, e.g. character

`read_excel(path, col_types = "text")`

Set each column individually

`read_excel(path,`

`col_types = c("text", "guess", "guess", "numeric")`

COLUMN TYPES

1. Get a vector of sheet names from the file path.

2. Set the vector names to be the sheet names.

3. Use `purrr::map_dfr()` to read multiple files into one data frame.

`skip`

`logical`

`date`

`list`

`TRUE`

`2`

`hello`

`1947-01-08`

`hello`

`FALSE`

`3.45`

`world`

`1956-10-21`

`1`

`s1`

`s1 s2`

`s3`

To read multiple sheets:

- Get a vector of sheet names from the file path.
- Set the vector names to be the sheet names.
- Use `purrr::map_dfr()` to read multiple files into one data frame.

write_sheet(data, ss = NULL, sheet = NULL)

Write a data frame into a new or existing Sheet.

gs4_create(name, ...)

`sheets = NULL`) Create a newSheet with a vector

of names, a data frame,

or a (named) list of data

frames.

`skip - " " or "_"`

`date - "D"`

`datetime - "T"`

`character - "C"`

`integer - "I"`

`double - "d"`

`numeric - "N"`

`list - "L"`

`cell - "C"` Returns

list of raw cell data.

For whole-file operations (e.g. renaming, sharing, placing within a folder), see the tidyverse

package `googledrive` at

googledrive.tidyverse.org

READ SHEETS

A	B	C	D	E
1	x1	x2	x3	x4
2	x	z	8	NA
3	y	7	9	10

x1	x2	x3	x4	x5
x	NA	z	8	NA
x	z	8	NA	9
y	7	NA	9	10

SHEETS METADATA

Column specifications define what data type each column of a file will be imported as.

Use the `col_types` argument of `read_excel()` to set the column specification.

Guess column types

To guess a column type `read_sheet()` looks at the first 1000 rows of data. Increase with the

`guess_max`.

`read_sheet(path, col_types = "C")`

Set all columns to same type, e.g. character

`read_sheet(path, guess_max = Inf)`

SPREADSHEET_ID

`https://docs.google.com/spreadsheets/d/SPREADSHEET_ID/edit#gid=SHEET_ID`

`gs4_get(ss)` Get spreadsheet meta data.

`gs4_find(...)` Get data on all spreadsheet files.

`sheet_properties(ss)` Get a tibble of properties for each worksheet. Also `sheet_names()`.

WRITE SHEETS

`write_sheet(data, ss = NULL, sheet = NULL)`

Write a data frame into a new or existing Sheet.

`gs4_create(name, ...)`

`sheets = NULL`) Create a newSheet with a vector

of names, a data frame, or a (named) list of data

frames.

`skip - " " or "_"`

`date - "D"`

`datetime - "T"`

`character - "C"`

`integer - "I"`

`double - "d"`

`numeric - "N"`

`list - "L"`

`cell - "C"` Returns

list of raw cell data.

For whole-file operations (e.g. renaming, sharing, placing within a folder), see the tidyverse

package `googledrive` at

googledrive.tidyverse.org

GOOGLESHEETS COLUMN SPECIFICATION

A	B	C	D	E
1	x1	x2	x3	x4
2	x	z	8	NA
3	y	7	9	10

x1	x2	x3	x4	x5
x	NA	z	8	NA
x	z	8	NA	9
y	7	NA	9	10

Guess column types

To guess a column type `read_sheet()` looks at the first 1000 rows of data. Increase with the

`guess_max`.

`read_sheet(path, col_types = "C")`

range_read()

Column specifications define what data type each column of a file will be imported as.

Use the `col_types` argument of `read_sheet()` to set the column specification.

`read_sheet(path, range = "B1:D2")`

Also use the `range` argument with cell specification functions

`cell_limits()`, `cell_rows()`, `cell_cols()`, and `anchored()`.

FILE LEVEL OPERATIONS

`googlesheets4` also offers ways to modify other aspects of Sheets (e.g. freeze rows, set column width, manage (work)sheets). Go to googlesheets4.tidyverse.org to read more.

For whole-file operations (e.g. renaming, sharing, placing within a folder), see the tidyverse package `googledrive` at googledrive.tidyverse.org.

READ EXCEL FILES



READ SHEETS



WRITE SHEETS



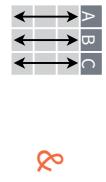
FILE LEVEL OPERATIONS



Data tidyng with `tidyverse`

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**



Reshape Data

- Pivot data to reorganize values into a new layout.

table4a		
country	year	cases
A	1999	0.7K
B	1999	3K
C	212K	213K

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")
```

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"))
```

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

```
fill(x, x2)
```

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 = NA) Add missing possible combinations of values of variables listed in ... Fill remaining variables with NA.

```
complete(mtcars, cyl, gear, carb)
```



Split Cells

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

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")
```



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.

Handle Missing Values

Drop or replace explicit missing values (NA).

na.rm = FALSE Collapse cells across several

columns into a single column.

na.omit(table5, century, year, col = "year", sep = "")

drop_na(data, ...) Drop columns.

replace_na(data, replace) Specify a value to replace NA in selected columns.

replace_na(x, list(x2 = 2))

Expand Tables



Nested Data

A **nested data frame** stores individual tables as a list-column of data frames within a larger organizing data frame. List-columns can also be lists of vectors or lists of varying data types.
Use a nested data frame to:

- Preserve relationships between observations and subsets of data. Preserve the type of the variables being nested (factors and datetimes aren't coerced to character).
- Manipulate many sub-tables at once with **purrr** functions like `map()`, `map2()`, or `pmap()` or with **dplyr** `rowwise()` grouping.

CREATE NESTED DATA

nest(`data, ...`) Moves groups of cells into a list-column of a data frame. Use alone or with `dplyr::group_by()`:

1. Group the data frame with `group_by()` and use `nest()` to move the groups into a list-column.

`n_storms <- storms %>% group_by(name) %>% nest()`

2. Use `nest(new_col = c(x, y))` to specify the columns to group using `dplyr::select()` syntax.

`n_storms <- storms %>% n_storms %>% select(name, films) %>% unnest_longer(films)`

`nest(data = c(year:long))`

`nest(data, cols, ..., keep_empty = FALSE)` Flatten nested columns back to regular columns. The inverse of `nest()`.

`n_storms %>% unnest(data)`

`unnest_longer(data, col, values_to = NULL, indices_to = NULL)` Turn each element of a list-column into a row.

`starwars %>% select(name, films) %>% unnest_wider(films)`

RESHAPE NESTED DATA

`unnest`(`data, cols, ...`, `keep_empty = FALSE`) Flatten nested columns back to regular columns. The inverse of `nest()`.

`n_storms %>% unnest(data)`

`unnest_longer`(`data, col, values_to = NULL, indices_to = NULL`) Turn each element of a list-column into a row.

TRANSFORM NESTED DATA

A vectorized function takes a vector, transforms each element in parallel, and returns a vector of the same length. By themselves Vectorized functions cannot work with lists, such as list-columns.

`dplyr::rowwise`(`data, ...`) Group data so that each row is one group, and within the groups, elements of list-columns appear directly (accessed with `[]`), not as lists of length one. When you **use `rowwise()`, `dplyr` functions will seem to apply functions to list-columns in a vectorized fashion.**

`unnest_wider`(`data, col`) Turn each element of a list-column into a regular column.

`starwars %>% select(name, films) %>% unnest_wider(films)`

`unnest`(`name, films`) %>%

`unnest_wider`(`films`)

`unnest`(`name`)

`unnest`(`films`)



APPLY A FUNCTION TO A LIST-COLUMN AND CREATE A REGULAR COLUMN

`n_storms %>% rowwise() %>% mutate(n = nrow(data))`

`mutate(n = nrow(data))`

`rowwise() %>% integer per row`

`append() returns a list for each row, so col type must be list`

`length() returns one integer per row`

`collapse multiple list-columns into a single list-column.`

`starwars %>% rowwise() %>% mutate(n_transports = length(c(vehicles, starships)))`

`apply a function to multiple list-columns.`

`starwars %>% rowwise() %>% mutate(transport = list(append(vehicles, starships)))`

`See purrr package for more list functions.`

CREATE TIBBLES WITH LIST-COLUMNS

tibble: tibble(...) Makes list-columns when needed.

`tibble(~max, ~seq, 3, 1:3, 4, 1:4, 5, 1:5)`

tibble: tibble(...) Saves list input as list-columns.

`tibble(max = c(3, 4, 5), seq = list(1:3, 1:4, 1:5))`

tibble: enframe(x, name = "name", value = "value")
Converts multi-level list to a tibble with list-cols.

`enframe(list('3'=1:3, '4'=1:4, '5'=1:5), 'max', 'seq')`

OUTPUT LIST-COLUMNS FROM OTHER FUNCTIONS

`dplyr::mutate(), transmute(), and summarise()` will output list-columns if they return a list.

`mtcars %>% group_by(cyl) %>% summarise(q = list(quartile(mpg)))`

`starwars %>% selectc(name, films) %>% hoist(first_film = 1, second_film = 2)`

`starwars %>% group_by(cyl) %>% summarise(q = list(quartile(mpg)))`

`Luke C-3PO R2-D2`

`Luke C-3PO R2-D2`

hoist(`data, col, ..., remove = TRUE`) Selectively pull list components out into their own top-level columns. Uses `purrr::pluck()` syntax for selecting from lists.

`starwars %>% selectc(name, films) %>% hoist(first_film = 1, second_film = 2)`

`starwars %>% group_by(cyl) %>% summarise(q = list(quartile(mpg)))`

`Luke C-3PO R2-D2`

`Luke C-3PO R2-D2`

collapse **multiple list-columns** into a single list-column.

`starwars %>% rowwise() %>% mutate(transport = list(append(vehicles, starships)))`

apply a function to multiple list-columns.

`starwars %>% rowwise() %>% mutate(transport = list(append(vehicles, starships)))`

`length() returns one integer per row`

`length() returns one integer per row`

Data transformation with dplyr :: CHEAT SHEET

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



Each **variable** is in **case**, is in its own **row** and **observation**, or its own **column**

Summarise Cases

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, ...)
summarise(mtcars, avg = mean(mpg))

count(.data, ..., wt = NULL, sort = FALSE, name = NULL) Count number of rows in each group defined by the variables in ... Also **tally()**.
count(mtcars, cyl)

Group Cases

Use **group_by**(.data, ...) add = FALSE, drop = TRUE) to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.

mtcars %>%
group_by(cyl) %>%
summarise(avg = mean(mpg))

ARRANGE CASES

Use **rowwise**(.data, ...) to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyverse cheat sheet for list-column workflow.

starwars %>%
rowwise() %>%
mutate(films_count = length(films))

ADD CASES

ungroup(x,...) Returns ungrouped copy of table.
ungroup(g_mtcars)

Manipulate Cases

EXTRACT CASES

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

filter(.data, ..., .preserve = FALSE) Extract rows that meet logical criteria.
filter(mtcars, mpg > 20)

distinct(.data, ..., .keep = all = FALSE) Remove rows with duplicate values.
distinct(mtcars, gear)

slice(.data, ..., .preserve = FALSE) Select rows by position.
slice(mtcars, 10:15)

slice_sample(.data, ..., n, prop, weight_by = NULL, replace = FALSE) Randomly select rows. Use n to select a number of rows and prop to select a fraction of rows.
slice_sample(mtcars, n = 5, replace = TRUE)

slice_min(.data, order_by = ..., n, prop, with_ties = TRUE) and **slice_max**() Select rows with the lowest and highest values.
slice_min(mtcars, mpg, prop = 0.25)

slice_head(.data, ..., n, prop) and **slice_tail**() Select the first or last rows.
slice_head(mtcars, n = 5)

LOGICAL AND BOOLEAN OPERATORS TO USE WITH FILTER()

== **<** **<=** **is.na()** **%in%** **|** **xor()**
!= **>** **>=** **!is.na()** **!** **&**

See ?base::Logic and ?Comparison for help.

Manipulate Variables

EXTRACT VARIABLES

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

pull(.data, var = -1, name = NULL, ...) Extract column values as a vector, by name or index.
pull(mtcars, wt)

select(.data, ...) Extract columns as a table.
select(mtcars, mpg, wt)

relocate(.data, ..., .before = NULL, .after = NULL) Move columns to new position.
relocate(mtcars, mpg, cyl, .after = last_col())

use_helpers(.data, .env = environment(.)) Use these helpers with **select()** and **across()**
e.g. **select(mtcars, mpg:cyl)** or **mutate(mtcars, across(everything(), mean))**

contains(match) **num_range**(prefix, range) ; e.g. mpg:cyl
ends_with(match) **all_of**(x) any_of(x, ..., vars) ; e.g. -gear
starts_with(match) **matches**(match) **everything()**

MANIPULATE MULTIPLE VARIABLES AT ONCE

across(cols, funs, ..., names = NULL) Summarise or mutate multiple columns in the same way.
summarise(mtcars, across(everything(), mean))

c_across(.cols) Compute across columns in row-wise manner.
transmute(rowwise(USgases), total = sum(c_across(1:2)))

MAKE NEW VARIABLES

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

VECTORIZED FUNCTION

arrange(.data, ..., by_group = FALSE) 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))

mutate(.data, ..., keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add_column()**, **add_count()**, and **add_tally()**.
mutate(mtcars, gpm = 1 / mpg)

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

rename(.data, ...) Rename columns. Use **rename_with()** to rename with a function.
rename(cars, distance = dist)



Vectorized Functions

Summary Functions

Combine Tables

COMBINE CASES

x	y
A B C	E F G
a t 1	a t 3
b u 2	= a t 1 a t 3
c v 3	b u 2 b u 2
d w 1	c v 3 d w 1

bind_rows(..., .id = NULL) Returns tables placed side by side as a single table. Column lengths must be equal. Columns will NOT be matched by id (to do that look at Relational Data below), so be sure to check that both tables are ordered the way you want before binding.



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.

COMBINE VARIABLES

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

bind_rows(..., .id = NULL) Returns tables placed side by side 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).

vectorized function

dplyr::lag() - offset elements by 1
dplyr::lead() - offset elements by -1

COUNT

dplyr::n() - number of values/rows
dplyr::n_distinct() - # of uniques
sum(is.na()) - # of non-NAs

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

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

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

anti_join(x, y, by = NULL, copy = FALSE, ...) **na_matches = "na"** Return rows of x that do not have a match in y. Use to see what what will not be included in a join.

semi_join(x, y, by = NULL, copy = FALSE, ...) **na_matches = "na"** Return rows of x that have a match in y. Use to see what will be included in a join.

POSITION

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

LOGICAL

cummax() - cumulative sum()
cumprod() - cumulative prod()
sum(is.na()) - # of non-NAs

nest_join(x, y, by = NULL, copy = FALSE, name = "name") Keep columns from both tables, nest them in a list, and then join them back together.

RELATIONAL DATA

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, na_matches = "na") Join matching values from y to x.

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

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

anti_join(x, y, by = NULL, copy = FALSE, ...) **na_matches = "na"** Return rows of x that do not have a match in y. Use to see what what will not be included in a join.

semi_join(x, y, by = NULL, copy = FALSE, ...) **na_matches = "na"** Return rows of x that have a match in y. Use to see what will be included in a join.

RANKING

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

RANK

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

dplyr::dense_rank() - rank w ties = min, no gaps

dplyr::min_rank() - rank with ties = min
dplyr::rank() - ranks 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

MISCELLANEOUS

dplyr::case_when() - multi-case if_else()

Row Names

starwars %>%
 mutate(type = case_when(
 height > 200 | mass > 200 ~ "large",
 species == "Droid" ~ "robot",
 ~ "other")
)

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()

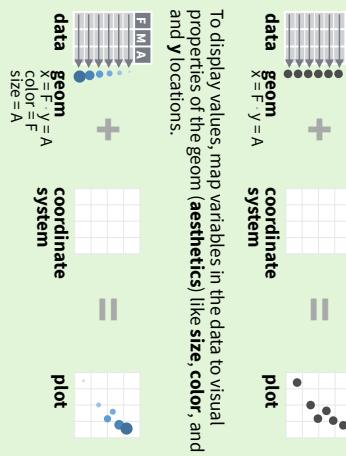
Also **tibble::column_torownames()**
Move col into row names.
tibble::column_torownames(a, var = "C")

tibble::column_torownames()
Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables.
left_join(x, y, by = c("C" = "D"))
suffix = c("1", "2"))

Data visualization With ggplot2 : CHEAT SHEET

Basics

ggplot2 is based on the **grammar of graphics**, the idea that you can build every graph from the same components: a **data set**, a **coordinate system**, and **geoms**—visual marks that represent data points.



To display values, map variables in the data to visual properties of the geom (**aesthetics**) like **size**, **color**, and **x** and **y** locations.

Complete the template below to build a graph.

Required

ggplot(**data** = <**DATA**>) +
<**GEOM_FUNCTION**>(**mapping** = **aes**(<**MAPPINGS**>,
stat = <**STAT**>, **position** = <**POSITION**>) +
<**COORDINATE_FUNCTION**> +
<**FACEIT_FUNCTION**> +
<**SCALE_FUNCTION**> +
<**THEME_FUNCTION**>

Geom

Use a geom function to represent data points, use the geom's aesthetic properties to represent variables. Each function returns a layer.

GRAPHICAL PRIMITIVES

a <- ggplot(economics, aes(date, unemploy))

b <- ggplot(mpg, aes(cty, hwy))

a + geom_blank() and a + expand_limits()

Ehnsure limits include values across all plots.

b + geom_curve(aes(yend = lat + 1,
xend = long + 1), curvature = 1) - x, xend, y, yend,
alpha, angle, color, curvature, linetype, size

a + geom_path(linewidth = "butt",
linejoin = "round") linemitre = 1)
x, y, alpha, color, group, linetype, size

a + geom_polygon(aes(alpha = 50)) - x, y, alpha,
color, fill, group, subgroup, linetype, size

b + geom_rect(aes(xmin = long, ymin = lat,
xmax = long + 1, ymax = lat + 1)) - xmax, xmin,
ymax, ymin, alpha, color, fill, linetype, size

a + geom_ribbon(aes(ymin = unemploy - 900,
ymax = unemploy + 900) - x, ymax, ymin,
alpha, color, fill, group, linetype, size

TWO VARIABLES both continuous

e <- ggplot(diamonds, aes(carat, price))

h <- ggplot(diamonds, aes(carat, price))

e + geom_label(aes(label = cty), nudge_x = 1,
nudge_y = 1) - x, y, label, alpha, angle, color,

B C family, fontface, hjust, lineheight, size, vjust

e + geom_point()

x, y, alpha, color, fill, shape, size, stroke

e + geom_quantile()

x, y, alpha, color, group, linetype, size, weight

e + geom_rug(sides = "bl")

x, y, alpha, color, linetype, size

e + geom_smooth(method = lm)

x, y, alpha, color, fill, group, linetype, size, weight

i + geom_area()

x, y, alpha, color, fill, linetype, size

i + geom_hex()

x, y, alpha, color, fill, size

continuous bivariate distribution

continuous function

i + geom_step(direction = "hv")

x, y, alpha, color, group, linetype, size

f + geom_col()

x, y, alpha, color, fill, group, linetype, size

f + geom_boxplot()

x, y, lower, middle, upper, ymax, ymin, alpha,

color, fill, group, linetype, shape, size, weight

f + geom_dodgeplot(binxaxis = "y", stackdir = "center")

x, y, alpha, color, fill, group

f + geom_violin(scale = "area")

x, y, alpha, color, fill, group, linetype, size, weight

both discrete

g <- ggplot(diamonds, aes(cut, color))

g + geom_count()

x, y, alpha, color, fill, shape, size, stroke

e + geom_jitter(height = 2, width = 2)

x, y, alpha, color, fill, shape, size

maps

data <- data.frame(murder = USArrests\$Murder,

- state = tolower(trownames(USArrests)))

map <- map_data("state")

k <- ggplot(data, aes(fill = murder))

k + geom_map(aes(map_id = state), map = map)

+ expand_limits(x = map\$long, y = map\$lat)

THREE VARIABLES

sealSz <- with(seals, sqrt(delta_long^2 + delta_lat^2)); l <- ggplot(seals, aes(long, lat))

l + geom_contour(aes(z = z))

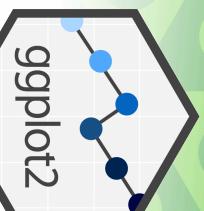
x, y, z, alpha, color, group, linetype, size, weight

vjust = 0.5, interpolate = FALSE)

x, y, alpha, fill

I + geom_tile(aes(fill = z))

x, y, alpha, color, fill, linetype, size, width



Stats

An alternative way to build a layer.

Scales

Override defaults with **scales** package.

Coordinate Systems

Faceting



A stat builds new variables to plot (e.g., count, prop).

data → **stat** → **geom** → **coordinate system**

geom to use → **stat function** → **geommapping**

geom = "polygon")
+ **stat** = density_2d(aes(fill = ..level..),
variable created by stat

n <- d + geom_bar()
Visualize a stat by changing the default stat of a geom
function, **geom_bar(stat = "count"**), or by using a stat
function, **stat_count(geom = "bar"**), which calls a default
geom to make a layer (equivalent to a geom function).
Use **..name..** syntax to map stat variables to aesthetics.

Scales map data values to the visual 'values' of an aesthetic. To change a mapping, add a new scale.

**n + scale_fill_manual(values = c("skyblue", "tealblue", "blue", "navy"),
limits = c("d", "e", "f", "g", "p", "r"),
name = "relin", labels = c("D", "E", "F", "G", "P", "R"))**
values to include in mapping
title to use in legend/ticks
labels to use in legend/ticks
breaks to use in legend/ticks

GENERAL PURPOSE SCALES

Use with most aesthetics

scale_*_continuous() - Map cont' values to visual ones.
scale_*_discrete() - Map discrete values to visual ones.

scale_*_binning() - Map continuous values to discrete bins.
scale_*_identity() - Use data values as visual ones.

c + stat_count(width = 1) **x, y** | ..count.., ..density..
c + stat_density(adjust = 1, kernel = "gaussian")
x, y | ..count.., ..density.., ..scaled..

e + stat_bin_2d(bins = 30, drop = T)
x, y, fill | ..count.., ..density..

e + stat_bin_hex(bins = 30) x, y, fill | ..count.., ..density..

l + stat_summary_2d(contour = TRUE, n = 100)
x, y, z, fill | ..value..

e + stat_summary_2d(aes(z = z), bins = 30, fun = mean)
x, y, z, fill | ..value..

f + stat_boxplot(coef = 1.5)
x, y | ..lower..., middle..., upper..., width..., ymin..., ymax..

f + stat_density(kernel = "gaussian", scale = "area") x, y
| ..density.., ..scaled.., ..count.., ..n.., ..id.., ..idlow.., ..width..

e + stat_ecdf(nn = 40) x, y | ..x..., ..y..

**e + stat_quantile(quantiles = c(0.1, 0.9),
formula = y ~ log(x), method = "iq") x, y | ..quantile..**

**e + stat_smooth(method = "lm", formula = y ~ x, se = T,
level = 0.95) x, y | ..se...., ..x..., ..ymin..., ..ymax..**

**ggplot() + xlim(-5, 5) + stat_function(fun = dnorm,
n = 20, geom = "point") x | ..x..., ..y..**

**ggplot() + stat_qq(aes(sample = 1:100))
x, y, sample | ..sample..., ..theoretical..**

**e + stat_summary(bin.fun = "mean", geom = "bar")
e + stat_identity()
e + stat_unique()**

SHAPE AND SIZE SCALES

p <- e + geom_point(aes(shape = fl, size = cyl))
+ **shape** = manual(values = c(3:7))
+ **size** = 4

**o + scale_shape_manual(values = 1:3)
o + scale_size_manual(values = 1:3)**

**o + scale_fill_gradient2(low = "red", high = "yellow",
mid = "white", midpoint = 25)**

**o + scale_fill_gradient(colors = topo.colors(6))
Also: rainbow(), heat.colors(), terrain.colors(), cm.colors(), RcolorBrewer::brewer.pal()**

r + theme_bw()
white background with grid lines.

r + theme_gray()
Grey background (default theme).

r + theme_light()

r + theme_minimal()
Minimal theme.

r + theme_void()
Empty theme.

r + theme_dark()
Dark for contrast.

r + theme() Customize aspects of the theme such
as axis, legend, panel, and facet properties.

r + ggtitle("Title") + theme(plot.title.position = "plot")

**t + scale_fill_discrete(name = "Title",
labels = c("A", "B", "C", "D", "E"))**
Set legend title and labels with a scale function.

**t + scale_x_continuous(limits = c(0, 100)) +
scale_y_continuous(limits = c(0, 100))**

Coordinate Systems

Faceting



Facets divide a plot into subplots based on the values of one or more discrete variables.

r <- d + geom_bar()
The default cartesian coordinate system.

r + coord_fixed(ratio = 1/2)
Fixed aspect ratio between x and y units.

ggplot(np3, aes(y = fl)) + geom_bar()
Flip cartesian coordinates by switching X and Y aesthetic mappings.

r + coord_polar(theta = "x", direction = 1)
theta, start, direction - Polar coordinates.

r + coord_trans(x = "sqrt") **x, y, xlim, ylim**
Transformed cartesian coordinates. Set **xtrans** and **ytrans** to the name of a windows function.

r + coord_quickmap()
projection = "ortho", orientation = c(41, -74, 0)
Map projections from the **maptools** package (mercator (default), azequalarea, lagrange, etc.).

Position Adjustments

Position adjustments determine how to arrange geoms that would otherwise occupy the same space.

s <- ggplot(mpg, aes(fl, fill = drv))
s + geom_bar(position = "dodge")

s + geom_bar(position = "stack")
Stack elements side by side.

s + geom_bar(position = "fill")
Stack elements on top of one another, normalize height.

e + geom_point(position = "jitter")
Add random noise to X and Y position of each element to avoid overplotting.

e + geom_label(position = "nudge")
Nudge labels away from points.

s + geom_bar(position = "stack")
Stack elements on top of one another.

Each position adjustment can be recast as a function with manual **width** and **height** arguments:

s + geom_bar(position = position_dodge(width = 1))

Labeled and Legends

Use **labs()** to label the elements of your plot.

**t + labs(x = "New x axis label", y = "New y axis label",
title = "Add a title above the plot",
subtitle = "Add a subtitle below title",
caption = "Add a caption below plot",
alt = "Add alt text to the plot",
aes = "New **aes** legend title")**

t + annotation(geom = "text", x = 8, y = 9, label = "A")
Places a geom with manually selected aesthetics.

t + guides(x = guide_axis(in.dodge = 2)) Avoid crowded or overlapping labels with guide_axis(in.dodge or angle).

n + guides(fill = "none") Set legend type for each aesthetic: colorbar, legend, or none (no legend).

**n + theme(panel.border = element_rect(colour = "black",
width = 1), panel.background = element_rect(fill = "white"))**
Place legend at "bottom", "top", "left" or "right".

**n + scale_fill_discrete(name = "Title",
labels = c("A", "B", "C", "D", "E"))**
Set legend title and labels with a scale function.

Themes

r + theme_bw()
white background with grid lines.

r + theme_light()

r + theme_minimal()
Minimal theme.

r + theme_void()
Empty theme.

r + theme_dark()
Dark for contrast.

r + theme() Customize aspects of the theme such
as axis, legend, panel, and facet properties.

r + ggtitle("Title") + theme(plot.title.position = "plot")

t + coord_cartesian(xlim = c(0, 100), ylim = c(10, 20))
With clipping (removes unseen data points):



rmarkdown :: CHEAT SHEET

What is rmarkdown?

 **Rmd files**: Develop your code and ideas side-by-side in a single document. Run code as individual chunks or as an entire document.

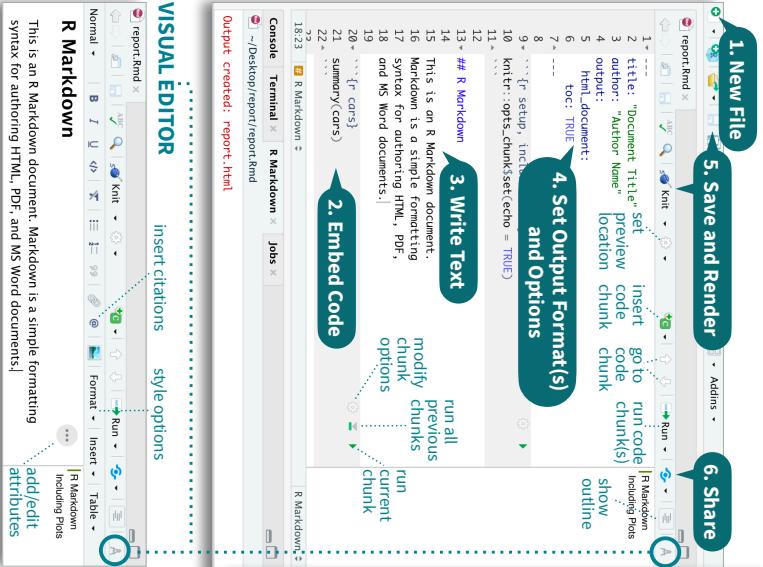
Dynamic Documents: Knit together plots, tables, and results with narrative text. Render to a variety of formats like HTML, PDF, MS Word, or MS Powerpoint.

Reproducible Research: Upload, link to, or attach your report to share. Anyone can read or run your code to reproduce your work.

Workflow

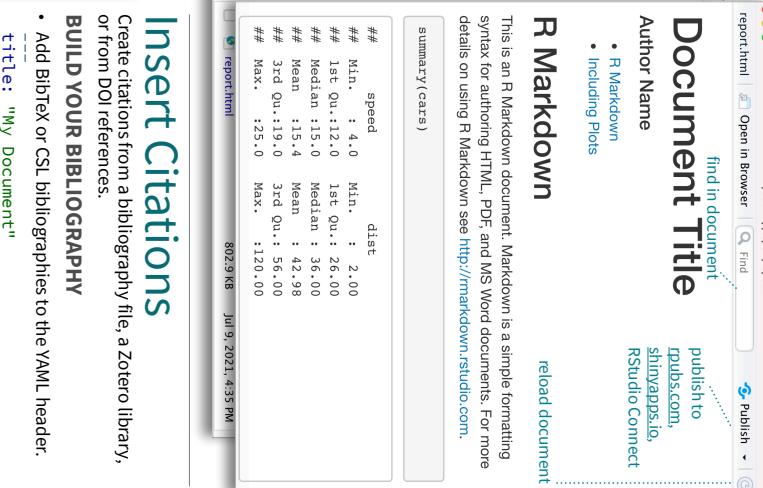
- 1 Open a new **Rmd file** in the RStudio IDE by going to **File > New File > R Markdown**.
- 2 **Embed code** in chunks. Run code by line, by chunk, or all at once.
- 3 Write **text** and add tables, figures, images, and citations. Format with Markdown syntax or the RStudio Visual Markdown Editor.
- 4 Set **output format(s)** and **options** in the YAML header. Customize themes or add parameters to execute or add interactivity with Shiny.
- 5 **Save and render** the whole document. Knit periodically to preview your work as you write.
- 6 **Share your work!**

SOURCE EDITOR



The screenshot shows the RStudio interface with the Source Editor tab selected. Step 1, "New File", is highlighted. Step 2, "Embed Code", is shown with a callout pointing to the code editor where a chunk of R code is being written. Step 3, "Write Text", is shown with a callout pointing to the text editor where "summary(cars)" is typed. Step 4, "Set Output Format(s) and Options", is shown with a callout pointing to the YAML header where output formats like "html_document" and "knitr::opts_chunk\$set(echo = TRUE)" are defined. Step 5, "Save and Render", is shown with a callout pointing to the "File" menu. Step 6, "Share", is shown with a callout pointing to the "Publish" button.

RENDERED OUTPUT



The screenshot shows the RStudio interface with the Rendered Output tab selected. It displays the generated HTML output for the R Markdown document. The output includes a summary of the cars dataset, such as "## # Min. : 4.0" and "## # Max. : 25.0". Below the summary, there are sections for "Header 1", "Header 2", "Header 3", "Header 6", "Header 6", "Header 6", and "Header 6". At the bottom, there is a section for "Caption" with a "markdown" icon.

Write with

Markdown

The syntax on the left renders as the output on the right.

Plain text.
End a line with two spaces to start a new paragraph.

Also end with a backslash \ to make a new line.

italics and **bold**
superscript²/subscript₂
~~strikethrough~~

italics and **bold**
superscript²/subscript₂
~~strikethrough~~
escaped: *`

endash: --, emdash: ---
#Header 1
#Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Header 1
Header 2
...
Header 6
- unnumbered list
- item 2
- item 2a (indent 1 tab)
- item 2b

Insert Citations

Create citations from a bibliography file, a Zotero library, or from DOI references.

BUILD YOUR BIBLIOGRAPHY

- Add BibTeX or CSL bibliographies to the YAML header.

```
title: "My Document"
bibliography: references.bib
link-citations: TRUE
```

- If Zotero is installed locally, your main library will automatically be available.

Add citations by searching "from DOI" in the **Insert Citation** dialog.

INSERT CITATIONS

- Access the **Insert Citations** dialog in the Visual Editor by clicking the @ symbol in the toolbar or by clicking **Insert > Citation**.
- Add citations with markdown syntax by typing **[@cite]** or **@cite**.

Set options for the entire document in the first chunk.

```
```{r include=FALSE}
knitr::opts_chunk$message = FALSE
```

```
```{r}
fig_align: "left", "right", or "center"
fig_collapse: NULL, "collapse"
fig_path: "figure/"
fig_width & fig_height: 7
out_width: FALSE
collapse comment: "#"
child: NULL
purp: TRUE
```

Insert `r **<code>**` into text sections. Code is evaluated at render and results appear as text.

"Built with **r getRVersion()**"

See more options and defaults by running **str(knitr::opts_chunk\$get()**)

Insert Tables

Output data frames as tables using **kable(data, caption)**.

```
```{r}
kable(data, caption = "Table with kable")
```

```
```{r}
data <- faithful[1:4, ]
kable(data, caption = "Table with kable")
```

```
```{r}
Results {dataset}
Plots text
text
Tables more text
```

## HTML Tables

### Results

### Plots

### Tables

### Text

### Other table packages include **flextable**, **gt**, and **kableExtra**.



# Set Output Formats and their Options in YAML

Use the document's YAML header to set an **output format** and customize it with **output options**.

	<b>IMPORTANT OPTIONS</b>	<b>DESCRIPTION</b>
<b>anchor_sections</b>	Show section anchors on mouse hover (TRUE or FALSE)	x
<b>citation_package</b>	The LaTeX package to process citations ("default", "natbib", "biblatex")	x
<b>code_download</b>	Let readers to toggle the display of R code ("none", "hide", or "show")	x
<b>css</b>	CSS or SCSS file to use to style document (e.g., "style.css")	x
<b>dev</b>	Graphics device to use for figure output (e.g., "png", "pdf")	x x
<b>df_print</b>	Method for printing data frames ("default", "kable", "tibble", "paged")	x x x x
<b>fig_caption</b>	Should figures be rendered with captions (TRUE or FALSE)	x x x x
<b>highlight</b>	Syntax highlighting ("tango", "pygments", "kate", "zenburn", "extmate")	x x x
<b>includes</b>	File of content to place in doc ("in_header", "before_body", "after_body")	x x
<b>keep_md</b>	Keep the Markdown .md file generated by knitting (TRUE or FALSE)	x x x
<b>keep_tex</b>	Keep the intermediate TEX file used to convert to PDF (TRUE or FALSE)	x x
<b>latex_engine</b>	LaTeX engine for producing PDF output ("pdflatex", "xelatex", or "ualatex")	x
<b>reference_docx_doc</b>	docx/pptx file containing styles to copy in the output (e.g., "file.docx", "file.pptx")	x x
<b>rstudio_presentation</b>	Theme options (see Bootswatch and Custom Themes below)	x
<b>slide_presentation</b>	Add a table of contents at start of document (TRUE or FALSE)	x x x x
<b>beamr_presentation</b>	The lowest level of headings to add to table of contents (e.g., 2, 3)	x x x x
<b>toc_depth</b>	Floating the table of contents to the left of the main document content (TRUE or FALSE)	x
<b>toc_float</b>	The lowest level of headings to add to table of contents (e.g., 2, 3)	x x x x

\* Requires LaTeX, use `tinytex::install_tinytex()`  
Also see `flexdashboard`, `bookdown`, `distrill`, and `blogdown`.

Use `?<output format>` to see all of a format's options, e.g. `?html_document`

## More Header Options

### PARAMETERS

Parameterize your documents to reuse with new inputs (e.g., data, values, etc.).

1. **Add parameters** ---  
`params:`  
  `state: "hawai'i"`  
  in the header as  
  sub-values of  
  params.
2. **Call parameters**  
`...{r}`  
`data <- df[, params$state]`  
`summary(data)`  
  in code using  
`params$name<-`

<b>CREATES</b>	<b>OUTPUT FORMAT</b>
.html	html_document
.pdf	pdf_document*
Microsoft Word (.docx)	word_document
Microsoft Powerpoint (.pptx)	powerpoint_presentation
OpenDocument Text	odt_document
Rich Text Format	rtf_document
Markdown	md_document
Markdown for Github	github_document
ioslides HTML slides	ioslides_presentation
Slidy HTML slides	slidy_presentation
Beamer slides	beamr_presentation

3. **Set parameters**  
with Knit with  
parameters or the  
params argument  
of render().



<b>REUSABLE TEMPLATES</b>
1. <b>Create a new package</b> with a <code>inst/rmarkdown/templates</code> directory.
2. <b>Add a folder</b> containing <code>template.yaml</code> (below) and <code>skel.Rmd</code> (template contents).

3. **Install** the package to access template by going to `File > New R Markdown > From Template`.

More on [bslib at pkgs.rstudio.com/bslib/](#).

# Render



When you render a document, `markdown`:

1. Runs the code and embeds results and text into an .md file with knitr.
2. Converts the .md file into the output format with Pandoc.



### HTML

HTML  
PDF  
MS Word  
MS PPT

### Save

Save, then **Knit** to preview the document output.

The resulting HTML/PDF/MS Word/etc. document will be created and saved in the same directory as the Rmd file.

### Publish

on **RStudio Connect** to share R Markdown documents securely, schedule automatic updates, and interact with parameters in real time.

### Share

This is a div with some text in it.

Also see [Shiny Prereendered](#) for better performance.

`markdowndown.rstudio.com/`

**authoring shiny\_prerendered**

Embed a complete app into your document with `shiny::shinyAppDir()`. More at [bookdown.org/yihui/markdown/shiny-embedded.html](#).

### BOOTSWATCH THEMES

Customize HTML documents with Bootswatch themes from the `bslib` package using the theme output option.

1. **Add parameters** ---  
`...{r}`  
`bslib::bootswatch_themes()` to list available themes.
2. Call `bslib::bootswatch_themes()` to get the theme name.

<b>CUSTOM THEMES</b>	<b>STYLING WITH CSS AND SCSS</b>
Customize individual HTML elements using bslib variables. Use <code>bslib_theme</code> to see more variables.	Add CSS and SCSS to your document by adding a path to a file with the <code>css</code> option in the YAML header.
1. Create a new package with a <code>inst/rmarkdown/templates</code> directory. 2. Add a folder containing <code>template.yaml</code> (below) and <code>skel.Rmd</code> (template contents). --- name: "My Template" fg: "#E4E4E4" base_font: google: "Prompt"	<p>1. Add <b>runtime: shiny</b> to the YAML header.</p> <pre>title: "My Document" author: "Author Name" output:   html_document:     css: "style.css"</pre> <p>2. Call Shiny input functions to embed input objects.</p> <p>3. Call Shiny render functions to embed reactive output.</p> <p>4. Render with <code>rmarkdown::run()</code> or click <b>Run Document</b> in RStudio IDE.</p>

### INTERACTIVITY

Turn your report into an interactive Shiny document in 4 steps:

1. Add **runtime: shiny** to the YAML header.
2. Call Shiny input functions to embed input objects.
3. Call Shiny render functions to embed reactive output.
4. Render with `rmarkdown::run()` or click **Run Document** in RStudio IDE.



### Bracketed Span

A [green]{my-color} word.



### Fenced Div

...{my-color}  
All of these words are green.



### How many cars?

5  
speed dist  
1 4.00 2.00  
2 4.00 10.00  
3 7.00 4.00  
4 7.00 22.00  
5 8.00 16.00



# String manipulation with stringr :: CHEAT SHEET

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

```
TRUE TRUE FALSE TRUE
↓ → ↓ → ↓
str_detect(string, pattern, negate = FALSE)
 Detect the presence of a pattern match in a
 string. Also str_like(). str_detect(fruit, "a")
```

```
TRUE TRUE FALSE TRUE
↓ → ↓ → ↓
str_starts(string, pattern, negate = FALSE)
 Detect the presence of a pattern match at
 the beginning of a string. Also str_ends().
 str_starts(fruit, "a")
```

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

```
start end
1 2 4 2 4
↓ → ↓ → ↓
str_locate(string, pattern)
 Locate the
 positions of pattern matches in a string.
 Also str_locate_all(). str_locate(fruit, "a")
```

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

## Mutate Strings

```
str_sub() <- value. Replace substrings by
 identifying the substrings with str_sub() and
 assigning into the results.
str_sub(fruit, 1, 3) <- "Str"
```

```
str_replace(string, pattern, replacement)
 Replace the first matched pattern in each
 string. Also str_remove().
str_replace(fruit, "p", "m")
```

```
str_replace_all(string, pattern, replacement)
 Replace all matched patterns in each string.
 Also str_remove_all().
str_replace_all(fruit, "p", "")
```

```
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)
```

```
a string
↓
A STRING
 A string
str_to_title(string, locale = "en")1
 Convert strings to title case. Also str_to_sentence().
str_to_title(sentences)
```

## Subset Strings

```
str_sub(string, start = 1L, end = -1L)
 Extract
 substrings from a character vector.
str_sub(fruit, 1, 3); str_sub(fruit, -2)
```

```
str_subset(string, pattern, negate = FALSE)
 Return only the strings that contain a pattern
 match. str_subset(fruit, "p")
```

```
str_extract(string, pattern)
 Return the first
 pattern match found in each string, as a vector.
 Also str_extract_all() to return every pattern
 match. str_extract(fruit, "[aeiou]")
```

```
str_match(string, pattern)
 Return the
 first pattern match found in each string, as
 a matrix with a column for each () group in
 pattern. Also str_match_all().
str_match(sentences, "(a|the) ([^])")
```

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

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

## Manage Lengths

```
str_length(string)
 The width of strings (i.e.
 the number of code points, which generally equals
 the number of characters). str_length(fruit)
```

```
str_pad(string, width, side = c("left", "right",
 "both"), pad = "")1
 Pad strings to constant
 width. str_pad(fruit, 17)
```

```
str_trunc(string, width, side = c("right", "left",
 "center"), ellipsis = "...")1
 Truncate the width
 of strings, replacing content with ellipsis.
str_trunc(sentences, 6)
```

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

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

## stringr

stringr  
version 1.0.7

stringr  
version 1.0.7

## Join and Split

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

```
str_flatten(string, collapse = "")1
 Combines
 into a single string, separated by collapse.
str_flatten(fruit, "")
```

```
str_dup(string, times)
 Repeat strings times
 times. Also str_unique() to remove duplicates.
str_dup(fruit, times = 2)
```

```
str_split_fixed(string, pattern, n)
 Split a
 vector of strings into a matrix of substrings
 (splitting at occurrences of a pattern match).
 Also str_split() to return a list of substrings
 and str_split_n() to return the nth substrings.
str_split_fixed(sentences, "", n=3)
```

```
appl <-->
banana
p <--> ar
 ↗
 ↘
```

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

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

```
str_equal(x, y, locale = "en", ignore_case =
 FALSE,...)1
 Determine if two strings are
 equivalent. str_equal(c("a", "b"), c("a", "C"))
 ↗
 ↘
```

```
str_wrap(string, width = 80, indent = 0,
 exdent = 0)1
 Wrap strings into nicely formatted
 paragraphs. str_wrap(sentences, 20)
```

## Order Strings

```
str_order(x, decreasing = FALSE, na_last =
 TRUE, locale = "en", numeric = FALSE,...)1
 Return the vector of indexes that sorts a
 character vector. fruit$str_order(fruit)
```

```
str_sort(x, decreasing = FALSE, na.last =
 TRUE, locale = "en", numeric = FALSE,...)1
 Sort a character vector. str_sort(fruit)
```

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

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

```
str_equal(x, y, locale = "en", ignore_case =
 FALSE,...)1
 Determine if two strings are
 equivalent. str_equal(c("a", "b"), c("a", "C"))
 ↗
 ↘
```

```
str_wrap(string, width = 80, indent = 0,
 exdent = 0)1
 Wrap strings into nicely formatted
 paragraphs. str_wrap(sentences, 20)
```

## Helpers

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

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

```
str_equal(x, y, locale = "en", ignore_case =
 FALSE,...)1
 Determine if two strings are
 equivalent. str_equal(c("a", "b"), c("a", "C"))
 ↗
 ↘
```

```
str_wrap(string, width = 80, indent = 0,
 exdent = 0)1
 Wrap strings into nicely formatted
 paragraphs. str_wrap(sentences, 20)
```

```
1 See bit.ly/ISO639-1 for a complete list of locales.
```

# Need to Know

Pattern arguments in `stringr` are interpreted as regular expressions *after any special characters have been parsed*.

In R, you write regular expressions as *strings*, sequences of characters surrounded by quotes ("") or single quotes('').

Some characters cannot be represented directly in an R string . These must be represented as **special characters**, sequences of characters that have a specific meaning, e.g.

Special Character	Represents
<code>\\"</code>	"
<code>\n</code>	new line
<code>\t</code>	tab
<code>\s</code>	any whitespace ( <code>\\$</code> for non-whitespace)
<code>\d</code>	any digit ( <code>\D</code> for non-digits)
<code>\w</code>	any word character ( <code>\W</code> for non-word chars)
<code>\b</code>	word boundaries
<code>[:digit:]</code>	digits
<code>[:alpha:]</code>	letters
<code>[:lower:]</code>	lowercase letters
<code>[:upper:]</code>	uppercase letters
<code>[:alnum:]</code>	letters and numbers
<code>[:punct:]</code>	punctuation
<code>[:graph:]</code>	letters, numbers, and punctuation
<code>[:space:]</code>	space characters (i.e. <code>\s</code> )
<code>[:blank:]</code>	space and tab (but not new line)
.	every character except a new line

## INTERPRETATION

Patterns in stringr are interpreted as regexes. 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 as end of strings, allow R comments within regex's, and/or have . match everything including `\n`.

`str_detect("", regex("\\", TRUE))`

**fixed()** Matches raw bytes but will miss some characters that can be represented in multiple ways (fast). `str_detect("\u0130", fixed("ö"))`

**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("ö", TRUE, locale = "tr"))`

**boundary()** Matches boundaries between characters, linebreaks, 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

`string` (type this)    `regexp` (to mean this)    `matches` (which matches this)

**a (etc.)**

a (etc.)

see "`\a`"

`\abc ABC 123 .!?\00`

**\b**

·

see "`\b`"

`\abc ABC 123 .!?\00`

**\!**

!

see "`\!`"

`\abc ABC 123 .!?\00`

**\?**

?

see "`\?`"

`\abc ABC 123 .!?\00`

**\|**

·

see "`\|\|`"

`\abc ABC 123 .!?\00`

**\{**

{

see "`\{\}`"

`\abc ABC 123 .!?\00`

**\}**

}

see "`\}\}`"

`\abc ABC 123 .!?\00`

**\\_**

\_

see "`\_\_`"

`\abc ABC 123 .!?\00`

**\n**

new line

see "`\n`"

`\abc ABC 123 .!?\00`

**\t**

tab

see "`\t`"

`\abc ABC 123 .!?\00`

**\s**

space

see "`\s`"

`\abc ABC 123 .!?\00`

**\d**

·

see "`\d`"

`\abc ABC 123 .!?\00`

**\w**

word

see "`\w`"

`\abc ABC 123 .!?\00`

**\b**

word boundaries

see "`\b`"

`\abc ABC 123 .!?\00`

**\d**

digits

see "`\d`"

`\abc ABC 123 .!?\00`

**\l**

letters

see "`\l`"

`\abc ABC 123 .!?\00`

**\L**

lowercase letters

see "`\L`"

`\abc ABC 123 .!?\00`

**\U**

uppercase letters

see "`\U`"

`\abc ABC 123 .!?\00`

**\A**

letters and numbers

see "`\A`"

`\abc ABC 123 .!?\00`

**\P**

punctuation

see "`\P`"

`\abc ABC 123 .!?\00`

**\G**

letters, numbers, and punctuation

see "`\G`"

`\abc ABC 123 .!?\00`

**\S**

space characters (i.e. `\s`)

see "`\S`"

`\abc ABC 123 .!?\00`

**\B**

space and tab (but not new line)

see "`\B`"

`\abc ABC 123 .!?\00`

**\Z**

every character except a new line

see "`\Z`"

`\abc ABC 123 .!?\00`

**alt <- function(rx, str\_view\_all("abcde", rx))**

example

`\abc \abcde`

`\abc ABC 123 .!?\00`

**QUANTIFIERS**

**quant <- function(rx, str\_view\_all(".aa.aaa", rx))**

example

`\aa \aaa \aaaa`

`\aa.\aa.\aaa \aa.\aa.\aaa \aa.\aa.\aaa`

**GROUPS**

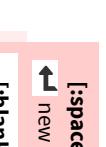
**ref <- function(rx, str\_view\_all("abbab", rx))**

example

`\ab\de \ab\de`

`\ab\de \ab\de \ab\de`

**stringr**



# Apply functions with purrr :: CHEAT SHEET

## Map Functions

ONE LIST

**map(x, f, ...)** Apply a function to each element of a list or vector; return a list.  
 $x \leftarrow \text{list}(1:10, 11:20, 21:30)$   
 $y \leftarrow \text{list}(1, 2, 3); z \leftarrow \text{list}(4, 5, 6); l2 \leftarrow \text{list}(x = "a", y = "Z")$   
 $\text{map}(l1, \text{sort}, \text{decreasing} = \text{TRUE})$



TWO LISTS

**map2(x, y, f, ...)** Apply a function to pairs of elements from two lists or vectors; return a list.  
 $y \leftarrow \text{list}(1, 2, 3); z \leftarrow \text{list}(4, 5, 6); l2 \leftarrow \text{list}(x = "a", y = "Z")$   
 $\text{map2}(x, y, \sim x * y)$



MANY LISTS

**pmap(l, f, ...)** Apply a function to groups of elements from a list of lists or vectors; return a list.  
 $\text{pmap}(\text{list}(x, y, z), \sim \dots 1 * (\dots 2 + \dots 3))$   
 $\text{imap}(y, \sim \text{paste0}(y, "\:", "\", x))$



LISTS AND INDEXES

**imap(x, f, ...)** Apply f to each element and its index; return a list.  
 $\text{imap}(y, \sim \text{paste0}(y, "\:", "\", x))$



## Function Shortcuts

Use `~.` with functions like `map()` that have single arguments.

`map(l, ~ . + 2)`

becomes

`map(l, function(x) x + 2))`

Use `~.` with functions like `map2()` that have two arguments.

`map2(l, p, ~ .x + y)`

becomes

`map2(l, p, function(x, y) x + y))`

Use `~.1..2..3` etc with functions like `pmap()` that have many arguments.

`pmap(list(a, b, c), ~ .3 + ..1 .. 2)`

becomes

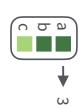
`pmap(list(a, b, c), function(a, b, c) c + a - b))`

Use a **string** or an **integer** with any map function to index list elements by name or position. `map(l, "name")` becomes `map(l, function(x) x[["name"]])`



# Work with Lists

## Filter

	<b>keep(x, p...)</b>	Select elements that pass a logical test.
	<b>Conversely, discard()</b>	Conversely, <b>discard()</b> .
	<b>compact(x, p = identity)</b>	Drop empty elements.
	<b>compact(x)</b>	Drop empty elements.
	<b>head_(x, p = ...)</b>	Return head elements until one does not pass.
	<b>tail_(x, p = ...)</b>	Also <b>tail_while()</b> .
	<b>head_while(x, is.character)</b>	head_while(x, is.character)

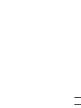
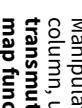
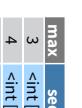
## Index

	<b>pluck_(x, ..., default=NULL)</b>	Select an element by name or index. Also <b>attr_getter()</b> and <b>chuck()</b> .
	<b>pluck_(x, "b")</b>	$x \%>%>% pluck("b")$
	<b>assign_(x, where, value)</b>	Assign a value to a location using <b>pluck</b> selection.
	<b>assign_in_(x, "b", 5)</b>	$x \%>%>% assign\_in("b", 5)$
	<b>modify_(x, .where, f...)</b>	Apply a function to a value at a selected location.
	<b>modify_in_(x, .where, f...)</b>	$x \%>%>% modify\_in("b", abs)$
	<b>modify_depth_(x, depth, f...)</b>	Apply function to each element at a given level of a list. Also <b>map_depth()</b>
	<b>modify_depth_(x, 2, ~+2)</b>	$modify\_depth(x, 2, ~+2)$

## Modify

	<b>modify_(x, f...)</b>	Apply a function to each element. Also <b>modify2()</b> , and <b>immodify()</b> .
	<b>modify2_(x, ~+2)</b>	
	<b>immodify_(x, ~+2)</b>	
		

## List-COLUMNS

	<b>flatten_(x)</b>	Remove a level of indexes from a list. Also <b>flatten_branch()</b> .
	<b>flatten_(x, margin = 3)</b>	Indexes from a list. Also <b>flatten_branch()</b> .
	<b>array_branch_(x, filter = NULL)</b>	All combinations of <b>x</b> and <b>y</b> . Also <b>cross()</b> , <b>cross3()</b> , and <b>cross_dfl()</b> .
	<b>cross2_(x, y, filter = NULL)</b>	All combinations of <b>x</b> and <b>y</b> . Also <b>cross()</b> , <b>cross3()</b> , and <b>cross_dfl()</b> .
	<b>transpose_(x, .names = NULL)</b>	Transposes the index ordering in a multi-level list.
	<b>transpose_(x)</b>	
	<b>set_names_(x, nm = x)</b>	Set the names of a vector/list directly or with a function.
	<b>vec_depth_(x)</b>	Return depth (number of levels of indexes).
	<b>set_names_(x, c("p", "q", "r"))</b>	Set the names of a vector/list directly or with a function.
	<b>set_names_(x, tolower)</b>	
	<b>accumulate_(x, f..., init)</b>	Reduce a list, but also return intermediate results. Also <b>accumulate2()</b> .
	<b>reduce_(x, f..., init, dir = c("forward", "backward"))</b>	Apply function recursively to each element of a list or vector. Also <b>reduce2()</b> .
	<b>reduce_(x, sum)</b>	



# Factors withforcats :: CHEAT SHEET

The **forcats** package provides tools for working with factors, which are R's data structure for categorical data.

## Factors

R represents categorical data with factors. A **factor** is an integer vector with a **levels** attribute that stores a set of mappings between integers and categorical values. When you view a factor, R displays not the integers, but the levels associated with them.

Create a factor with **factor()**

```
a = factor(c("a", "b", "c"))
a
#> [1] a b c
```

**fct\_relevel(f, ..., after=0L)**  
Manually reorder factor levels.

```
fct_relevel(f, c("b", "c", "a"))
fct_relevel(f, c("b", "c", "a"))
```

**fct\_infreq(f, ordered=NA)**  
Reorder levels by the frequency in which they appear in the data (highest frequency first).

```
f3 <- fct_inseq()
f3 <- fct_inseq(f3)
fct_infreq(f3)
```

**fct\_inorder(f, ordered=NA)**  
Reorder levels by order in which they appear in the data.

```
fct_inorder(f2)
f4 <- factor(c("a", "b", "c"))
fct_rev(f4)
```

## Inspect Factors

```
fct_levels(x) Return/set the levels of a factor. (levels(f); levels(f) <- c("x", "y", "z"))
Use unclass() to see its structure
```

```
a = factor(c("a", "b", "c"))
a
#> [1] a b c
```

```
b = factor(c("a", "b", "c"))
b
#> [1] a b c
```

```
c = factor(c("a", "b", "c"))
c
#> [1] a b c
```

**fct\_count(f, sort=FALSE, prop=FALSE)** Count the number of values with each level. **fct\_count(f)**

```
fct_match(f, lvls) Check for lvls in f. fct_match(f, "a")
```

```
fct_unique(f) Return the unique values, removing duplicates. fct_unique(f)
```

## Combine Factors

```
fct_c(...) Combine factors with different levels. Also fct_cross().
```

```
f1 <- factor(c("a", "c"))
f2 <- factor(c("b", "a"))
fct_c(f1, f2)
```

```
fct_unify(fs, levels=lvls, union(fs)) Standardize levels across a list of factors. fct_unify(lis(f2, f1))
```

```
a = factor(c("a", "b", "c"))
a
#> [1] a b c
```

```
b = factor(c("a", "b", "c"))
b
#> [1] a b c
```

```
c = factor(c("a", "b", "c"))
c
#> [1] a b c
```

```
a = fct_c(a, b)
a
#> [1] a b c
```

```
b = fct_c(b, c)
b
#> [1] a b c
```

```
a = fct_c(a, b)
a
#> [1] a b c
```

```
b = fct_c(b, c)
b
#> [1] a b c
```

```
c = fct_c(c, a)
c
#> [1] a b c
```

## Change the value of levels

**fct\_recode(f, ...)** Manually change levels. Also **fct\_relabel()** which obeys purrr::map syntax to apply a function or expression to each level.

```
fct_recode(f, v = "a", x = "b", z = "c")
fct_relabel(f, ~ paste0("X", .x))
```

**fct\_anon(f, prefix = "")**  
Anonymize levels with random integers.

```
fct_anon(f)
fct_anon(f)
```

**fct\_lump\_min(f, min, w = NULL, other\_level = "Other")** Lumps together factors that appear fewer than min times. Also **fct\_lump\_n0()**, **fct\_lump\_prop()**, and **fct\_lump\_lowfreq()**.

```
fct_lump_min(f, min = 2)
```

**fct\_lump\_min(f, min, w = NULL, other\_level = "Other")** Lumps together factors that appear fewer than min times. Also **fct\_lump\_n0()**, **fct\_lump\_prop()**, and **fct\_lump\_lowfreq()**.

```
fct_lump_min(f, min = 2)
```

**fct\_other(f, keep, drop, other\_level = "Other")** Replace levels with "other": **fct\_other(f, keep = c("a", "b"))**

```
fct_other(f, keep = c("a", "b"))
```

**fct\_reorder(f, x, fun = median, desc = FALSE)** Reorder levels by their relationship with another variable.

```
boxplot(data = PlantGrowth,
 weight ~ reorder(group, weight))
```

**fct\_reorder2(f, x, y, fun = last2, ..., desc = TRUE)** Reorder levels by their final values when plotted with two other variables.

```
ggplot(diamonds,aes(carat, price,
 color = fct_reorder2(color, carat,
 price))) + geom_smooth()
```

## Add or drop levels

**fct\_drop(f, only)** Drop unused levels.

```
f5 <- factor(c("a", "b", "x"))
f6 <- fct_drop(f5)
```

**fct\_expand(f, ..., x)** Add levels to a factor.

```
fct_expand(f6, x)
```

**fct\_explicit\_na(f, na\_level = "(Missing)")** Assigns a level to NAs to ensure they appear in plots, etc.

```
fct_explicit_na(factor(c("a", "b", NA)))
```

