

Data import with the tidyverse :: CHEATSHEET



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
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
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	1,5;2;3	4,5;5;NA
1	2	3
4	5	NA

read_csv2("file2.csv") Read semicolon delimited files with comma decimal marks.

`write_file("A;B;C\n1,5;2;3\n4,5;5;NA", file = "file2.csv")`

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

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

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")`

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
- **DBI** - databases
- **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.

```
spec(x)  
# cols(  
#   age = col_integer(),  
#   edu = col_character(),  
#   earn = col_double()  
# )
```

age is an integer

edu is a character

earn is a double (numeric)

COLUMN TYPES

Each column type has a function and corresponding string abbreviation.

- **col_logical()** - "l"
- **col_integer()** - "i"
- **col_double()** - "d"
- **col_number()** - "n"
- **col_character()** - "c"
- **col_factor(levels, ordered = FALSE)** - "f"
- **col_datetime(format = "")** - "T"
- **col_date(format = "")** - "D"
- **col_time(format = "")** - "t"
- **col_skip()** - "-", "_"
- **col_guess()** - "?"

USEFUL COLUMN ARGUMENTS

Hide col spec message

`read_*(file, show_col_types = FALSE)`

Select columns to import

Use names, position, or selection helpers.
`read_*(file, col_select = c(age, earn))`

Guess column types

To guess a column type, `read_*`() looks at the first 1000 rows of data. Increase with **guess_max**.
`read_*(file, guess_max = Inf)`

DEFINE COLUMN SPECIFICATION

Set a default type

```
read_csv(  
  file,  
  col_type = list(default = col_double()))  
)
```

Use column type or string abbreviation

```
read_csv(  
  file,  
  col_type = list(x = col_double(), y = "l", z = "_"))  
)
```

Use a single string of abbreviations

```
# col types: skip, guess, integer, logical, character  
read_csv(  
  file,  
  col_type = "_?ilc")  
)
```

Import Spreadsheets with readxl

READ EXCEL FILES

	A	B	C	D	E
1	x1	x2	x3	x4	x5
2	x		z	8	
3	y	7		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.xlsx")
```

READ SHEETS

A	B	C	D	E
s1	s2	s3		

s1	s2	s3
----	----	----

A	B	C	D	E
A	B	C	D	E

s1	s2	s3
----	----	----

```
path <- "your_file_path.xlsx"  
path >  
  excel_sheets() >  
  set_names() >  
  map(read_excel, path = path) >  
  list_rbind()
```

OTHER USEFUL EXCEL PACKAGES

For functions to write data to Excel files, see:

- **openxlsx**
- **writexl**

For working with non-tabular Excel data, see:

- **tidyxl**



with googlesheets4

READ SHEETS

	A	B	C	D	E
1	x1	x2	x3	x4	x5
2	x		z	8	
3	y	7		9	10

```
read_sheet(ss, sheet = NULL, range = NULL)  
Read a sheet from a URL, a Sheet ID, or a dribble  
from the googledrive package. See front page for  
more read arguments. Same as range_read().
```

SHEETS METADATA

URLs are in the form:
<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

1	x	4
2	y	5
3	z	6

A	B	C
1	1	x
2	2	y
3	3	z

1	A	B	C
2			

x1	x2	x3
2	y	5
3	z	6

A	B	C
1	x1	x2
2	1	x
3	2	y
4	3	z

write_sheet(data, ss = NULL, sheet = NULL)
Write a data frame into a new or existing Sheet.

gs4_create(name, ..., sheets = NULL) Create a new Sheet with a vector of names, a data frame, or a (named) list of data frames.

sheet_append(ss, data, sheet = 1) Add rows to the end of a worksheet.



GOOGLESHEETS4 COLUMN SPECIFICATION

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

Use the **col_types** argument of **read_sheet()**/ **range_read()** to set the column specification.

Guess column types

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_sheet(path, col_types = "text")`

Set each column individually

col types: skip, guess, integer, logical, character
`read_sheets(ss, col_types = "?ilc")`

COLUMN TYPES

logical	numeric	text	date	list
TRUE	2	hello	1947-01-08	hello
FALSE	3.45	world	1956-10-21	1

- skip
- guess
- logical
- numeric
- date
- list
- text

Use **list** for columns that include multiple data types. See **tidy** and **purrr** for list-column data.

1	x	4
2	y	5
3	z	6

1	A	B	C
2			

x1	x2	x3
2	y	5
3	z	6

A	B	C
1	x1	x2
2	1	x
3	2	y
4	3	z

write_sheet(data, ss = NULL, sheet = NULL)
Write a data frame into a new or existing Sheet.

gs4_create(name, ..., sheets = NULL) Create a new Sheet with a vector of names, a data frame, or a (named) list of data frames.

sheet_append(ss, data, sheet = 1) Add rows to the end of a worksheet.

- skip - "_" or "-"
- guess - "?"
- logical - "l"
- integer - "i"
- double - "d"
- numeric - "n"
- date - "D"
- datetime - "T"
- character - "c"
- list-column - "L"
- cell - "C" Returns list of raw cell data.

Use list for columns that include multiple data types. See **tidy** and **purrr** for list-column data.

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.

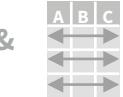
Data transformation with dplyr :: CHEATSHEET



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 $f(x, y)$

Summarize 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

summarize(.data, ...)
Compute table of summaries.
mtcars |> summarize(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()**, **add_count()**, **add_tally()**.
mtcars |> count(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) |> summarize(avg = mean(mpg))

Alternate grouping syntax with **.by** as an argument:

mtcars |> summarize(avg = mean(mpg), .by = cyl)

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

starwars |> rowwise() |> mutate(film_count = length(films))

ungroup(x, ...) Returns ungrouped copy of table.
g_mtcars <- mtcars |> group_by(cyl)
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.
mtcars |> filter(mpg > 20)



distinct(.data, ..., .keep_all = FALSE) Remove rows with duplicate values.
mtcars |> distinct(gear)



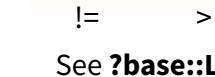
slice(.data, ..., .preserve = FALSE) Select rows by position.
mtcars |> slice(10:15)



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.
mtcars |> slice_sample(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.
mtcars |> slice_min(mpg, prop = 0.25)



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

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



add_row(.data, ..., .before = NULL, .after = NULL)
Add one or more rows to a table.
cars |> add_row(speed = 1, dist = 1)

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.
mtcars |> pull(wt)



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



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

Use these helpers with select() and across()

e.g. mtcars |> select(mpg:cyl)

contains(match)
ends_with(match)
starts_with(match)

num_range(prefix, range)
all_of(x)/any_of(x, ..., vars)
matches(match)

; e.g., mpg:cyl
!, e.g., !gear
everything()

MANIPULATE MULTIPLE VARIABLES AT ONCE

df <- tibble(x_1 = c(1, 2), x_2 = c(3, 4), y = c(4, 5))



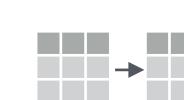
across(.cols, .funs, ..., .names = NULL) Summarize or mutate multiple columns in the same way.
df |> summarize(across(everything(), mean))



c_across(.cols) Compute across columns in row-wise data.
df |> rowwise() |> mutate(x_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
mutate(.data, ..., .keep = "all", .before = NULL, .after = NULL) Compute new column(s). Also **add_column()**.
mtcars |> mutate(gpm = 1 / mpg)
mtcars |> mutate(gpm = 1 / mpg, .keep = "none")



rename(.data, ...) Rename columns. Use **rename_with()** to rename with a function.
mtcars |> rename(miles_per_gallon = mpg)



Vectorized Functions

TO USE WITH MUTATE ()

mutate() applies 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

OFFSET

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

CUMULATIVE AGGREGATE

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

RANKING

dplyr::cume_dist() - proportion of all values <= 1
dplyr::dense_rank() - rank w 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

MISCELLANEOUS

dplyr::case_when() - multi-case if_else()
starwars |>
mutate(type = case_when(
height > 200 | mass > 200 ~ "large",
species == "Droid" ~ "robot",
TRUE ~ "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()

Summary Functions

TO USE WITH SUMMARIZE ()

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

summary function

COUNT

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

POSITION

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

LOGICAL

mean() - proportion of TRUEs
sum() - # of TRUEs

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.

tibble::rownames_to_column()
Move row names into col.
a <- mtcars |>
rownames_to_column(var = "C")

tibble::column_to_rownames()
Move col into row names.
a |> column_to_rownames(var = "C")

Also **tibble::has_rownames()** and **tibble::remove_rownames()**.

Combine Tables

COMBINE VARIABLES

X + y = <img alt="Diagram showing two tables X and y being combined into a single table with columns A, B, C, E, F, G. The first row of X (a t 1) is paired with the first row of y (a t 3). The second row of X (b u 2) is paired with the second row of y (b u 2). The third row of X (c v 3) is paired with the third row of y (d w 1). The fourth row of X (d w 1) is paired with the fourth row of y (a t 3). The fifth row of X (a t 3) is paired with the fifth row of y (b u 2). The sixth row of X (b u 2) is paired with the sixth row of y (c v 3). The seventh row of X (c v 3) is paired with the seventh row of y (d w 1). The eighth row of X (d w 1) is paired with the eighth row of y (a t 3). The ninth row of X (a t 3) is paired with the ninth row of y (b u 2). The tenth row of X (b u 2) is paired with the tenth row of y (c v 3). The eleventh row of X (c v 3) is paired with the eleventh row of y (d w 1). The twelfth row of X (d w 1) is paired with the twelfth row of y (a t 3). The thirteenth row of X (a t 3) is paired with the thirteenth row of y (b u 2). The fourteenth row of X (b u 2) is paired with the fourteenth row of y (c v 3). The fifteenth row of X (c v 3) is paired with the fifteenth row of y (d w 1). The sixteenth row of X (d w 1) is paired with the sixteenth row of y (a t 3). The seventeenth row of X (a t 3) is paired with the seventeenth row of y (b u 2). The eighteenth row of X (b u 2) is paired with the eighteenth row of y (c v 3). The nineteenth row of X (c v 3) is paired with the nineteenth row of y (d w 1). The twentieth row of X (d w 1) is paired with the twentieth row of y (a t 3). The twenty-first row of X (a t 3) is paired with the twenty-first row of y (b u 2). The twenty-second row of X (b u 2) is paired with the twenty-second row of y (c v 3). The twenty-third row of X (c v 3) is paired with the twenty-third row of y (d w 1). The twenty-fourth row of X (d w 1) is paired with the twenty-fourth row of y (a t 3). The twenty-fifth row of X (a t 3) is paired with the twenty-fifth row of y (b u 2). The twenty-sixth row of X (b u 2) is paired with the twenty-sixth row of y (c v 3). The twenty-seventh row of X (c v 3) is paired with the twenty-seventh row of y (d w 1). The twenty-eighth row of X (d w 1) is paired with the twenty-eighth row of y (a t 3). The twenty-ninth row of X (a t 3) is paired with the twenty-ninth row of y (b u 2). The thirtieth row of X (b u 2) is paired with the thirtieth row of y (c v 3). The thirtieth row of X (c v 3) is paired with the thirtieth row of y (d w 1). The thirty-first row of X (d w 1) is paired with the thirty-first row of y (a t 3). The thirty-second row of X (a t 3) is paired with the thirty-second row of y (b u 2). The thirty-third row of X (b u 2) is paired with the thirty-third row of y (c v 3). The thirty-fourth row of X (c v 3) is paired with the thirty-fourth row of y (d w 1). The thirty-fifth row of X (d w 1) is paired with the thirty-fifth row of y (a t 3). The thirty-sixth row of X (a t 3) is paired with the thirty-sixth row of y (b u 2). The thirty-seventh row of X (b u 2) is paired with the thirty-seventh row of y (c v 3). The thirty-eighth row of X (c v 3) is paired with the thirty-eighth row of y (d w 1). The thirty-ninth row of X (d w 1) is paired with the thirty-ninth row of y (a t 3). The forty-first row of X (a t 3) is paired with the forty-first row of y (b u 2). The forty-second row of X (b u 2) is paired with the forty-second row of y (c v 3). The forty-third row of X (c v 3) is paired with the forty-third row of y (d w 1). The forty-fourth row of X (d w 1) is paired with the forty-fourth row of y (a t 3). The forty-fifth row of X (a t 3) is paired with the forty-fifth row of y (b u 2). The forty-sixth row of X (b u 2) is paired with the forty-sixth row of y (c v 3). The forty-seventh row of X (c v 3) is paired with the forty-seventh row of y (d w 1). The forty-eighth row of X (d w 1) is paired with the forty-eighth row of y (a t 3). The forty-ninth row of X (a t 3) is paired with the forty-ninth row of y (b u 2). The五十th row of X (b u 2) is paired with the五十th row of y (c v 3). The fifty-first row of X (c v 3) is paired with the fifty-first row of y (d w 1). The fifty-second row of X (d w 1) is paired with the fifty-second row of y (a t 3). The fifty-third row of X (a t 3) is paired with the fifty-third row of y (b u 2). The fifty-fourth row of X (b u 2) is paired with the fifty-fourth row of y (c v 3). The fifty-fifth row of X (c v 3) is paired with the fifty-fifth row of y (d w 1). The fifty-sixth row of X (d w 1) is paired with the fifty-sixth row of y (a t 3). The fifty-seventh row of X (a t 3) is paired with the fifty-seventh row of y (b u 2). The fifty-eighth row of X (b u 2) is paired with the fifty-eighth row of y (c v 3). The fifty-ninth row of X (c v 3) is paired with the fifty-ninth row of y (d w 1). The六十th row of X (d w 1) is paired with the六十th row of y (a t 3). The六十-first row of X (a t 3) is paired with the六十-first row of y (b u 2). The六十-second row of X (b u 2) is paired with the六十-second row of y (c v 3). The六十-third row of X (c v 3) is paired with the六十-third row of y (d w 1). The六十-fourth row of X (d w 1) is paired with the六十-fourth row of y (a t 3). The六十-fifth row of X (a t 3) is paired with the六十-fifth row of y (b u 2). The六十-sixth row of X (b u 2) is paired with the六十-sixth row of y (c v 3). The六十第七 row of X (c v 3) is paired with the六十第七 row of y (d w 1). The六十第八 row of X (d w 1) is paired with the六十第八 row of y (a t 3). The六十第九 row of X (a t 3) is paired with the六十第九 row of y (b u 2). The六十第十 row of X (b u 2) is paired with the六十第十 row of y (c v 3). The六十第十一 row of X (c v 3) is paired with the六十第十一 row of y (d w 1). The六十第十二 row of X (d w 1) is paired with the六十第十二 row of y (a t 3). The六十第十三 row of X (a t 3) is paired with the六十第十三 row of y (b u 2). The六十第十四 row of X (b u 2) is paired with the六十第十四 row of y (c v 3). The六十第十五 row of X (c v 3) is paired with the六十第十五 row of y (d w 1). The六十第十六 row of X (d w 1) is paired with the六十第十六 row of y (a t 3). The六十第十七 row of X (a t 3) is paired with the六十第十七 row of y (b u 2). The六十第十八 row of X (b u 2) is paired with the六十第十八 row of y (c v 3). The六十第十九 row of X (c v 3) is paired with the六十第十九 row of y (d w 1). The六十第二十 row of X (d w 1) is paired with the六十第二十 row of y (a t 3). The六十第二十一 row of X (a t 3) is paired with the六十第二十一 row of y (b u 2). The六十第二十二 row of X (b u 2) is paired with the六十第二十二 row of y (c v 3). The六十第二十三 row of X (c v 3) is paired with the六十第二十三 row of y (d w 1). The六十第二十四 row of X (d w 1) is paired with the六十第二十四 row of y (a t 3). The六十第二十五 row of X (a t 3) is paired with the六十第二十五 row of y (b u 2). The六十第二十六 row of X (b u 2) is paired with the六十第二十六 row of y (c v 3). The六十第二十七 row of X (c v 3) is paired with the六十第二十七 row of y (d w 1). The六十第二十八 row of X (d w 1) is paired with the六十第二十八 row of y (a t 3). The六十第二十九 row of X (a t 3) is paired with the六十第二十九 row of y (b u 2). The六十第三十 row of X (b u 2) is paired with the六十第三十 row of y (c v 3). The六十第三十一 row of X (c v 3) is paired with the六十第三十一 row of y (d w 1). The六十第三十二 row of X (d w 1) is paired with the六十第三十二 row of y (a t 3). The六十第三十三 row of X (a t 3) is paired with the六十第三十三 row of y (b u 2). The六十第三十四 row of X (b u 2) is paired with the六十第三十四 row of y (c v 3). The六十第三十五 row of X (c v 3) is paired with the六十第三十五 row of y (d w 1). The六十第三十六 row of X (d w 1) is paired with the六十第三十六 row of y (a t 3). The六十第三十七 row of X (a t 3) is paired with the六十第三十七 row of y (b u 2). The六十第三十八 row of X (b u 2) is paired with the六十第三十八 row of y (c v 3). The六十第三十九 row of X (c v 3) is paired with the六十第三十九 row of y (d w 1). The六十第四十 row of X (d w 1) is paired with the六十第四十 row of y (a t 3). The六十第四十一 row of X (a t 3) is paired with the六十第四十一 row of y (b u 2). The六十第四十二 row of X (b u 2) is paired with the六十第四十二 row of y (c v 3). The六十第四十三 row of X (c v 3) is paired with the六十第四十三 row of y (d w 1). The六十第四十四 row of X (d w 1) is paired with the六十第四十四 row of y (a t 3). The六十第四十五 row of X (a t 3) is paired with the六十第四十五 row of y (b u 2). The六十第四十六 row of X (b u 2) is paired with the六十第四十六 row of y (c v 3). The六十第四十七 row of X (c v 3) is paired with the六十第四十七 row of y (d w 1). The六十第四十八 row of X (d w 1) is paired with the六十第四十八 row of y (a t 3). The六十第四十九 row of X (a t 3) is paired with the六十第四十九 row of y (b u 2). The六十第五十 row of X (b u 2) is paired with the六十第五十 row of y (c v 3). The六十第五十一 row of X (c v 3) is paired with the六十第五十一 row of y (d w 1). The六十第五十二 row of X (d w 1) is paired with the六十第五十二 row of y (a t 3). The六十第五十三 row of X (a t 3) is paired with the六十第五十三 row of y (b u 2). The六十第五十四 row of X (b u 2) is paired with the六十第五十四 row of y (c v 3). The六十第五十五 row of X (c v 3) is paired with the六十第五十五 row of y (d w 1). The六十第五十六 row of X (d w 1) is paired with the六十第五十六 row of y (a t 3). The六十第五十七 row of X (a t 3) is paired with the六十第五十七 row of y (b u 2). The六十第五十八 row of X (b u 2) is paired with the六十第五十八 row of y (c v 3). The六十第五十九 row of X (c v 3) is paired with the六十第五十九 row of y (d w 1). The六十第六十 row of X (d w 1) is paired with the六十第六十 row of y (a t 3). The六十第六十一 row of X (a t 3) is paired with the六十第六十一 row of y (b u 2). The六十第六十二 row of X (b u 2) is paired with the六十第六十二 row of y (c v 3). The六十第六十三 row of X (c v 3) is paired with the六十第六十三 row of y (d w 1). The六十第六十四 row of X (d w 1) is paired with the六十第六十四 row of y (a t 3). The六十第六十五 row of X (a t 3) is paired with the六十第六十五 row of y (b u 2). The六十第六十六 row of X (b u 2) is paired with the六十第六十六 row of y (c v 3). The六十第六十七 row of X (c v 3) is paired with the六十第六十七 row of y (d w 1). The六十第六十八 row of X (d w 1) is paired with the六十第六十八 row of y (a t 3). The六十第六十九 row of X (a t 3) is paired with the六十第六十九 row of y (b u 2). The六十第七十 row of X (b u 2) is paired with the六十第七十 row of y (c v 3). The六十第七十一 row of X (c v 3) is paired with the六十第七十一 row of y (d w 1). The六十第七十二 row of X (d w 1) is paired with the六十第七十二 row of y (a t 3). The六十第七十三 row of X (a t 3) is paired with the六十第七十三 row of y (b u 2). The六十第七十四 row of X (b u 2) is paired with the六十第七十四 row of y (c v 3). The六十第七十五 row of X (c v 3) is paired with the六十第七十五 row of y (d w 1). The六十第七十六 row of X (d w 1) is paired with the六十第七十六 row of y (a t 3). The六十第七十七 row of X (a t 3) is paired with the六十第七十七 row of y (b u 2). The六十第七十八 row of X (b u 2) is paired with the六十第七十八 row of y (c v 3). The六十第七十九 row of X (c v 3) is paired with the六十第七十九 row of y (d w 1). The六十第八十 row of X (d w 1) is paired with the六十第八十 row of y (a t 3). The六十第八十一 row of X (a t 3) is paired with the六十第八十一 row of y (b u 2). The六十第八十二 row of X (b u 2) is paired with the六十第八十二 row of y (c v 3). The六十第八十三 row of X (c v 3) is paired with the六十第八十三 row of y (d w 1). The六十第八十四 row of X (d w 1) is paired with the六十第八十四 row of y (a t 3). The六十第八十五 row of X (a t 3) is paired with the六十第八十五 row of y (b u 2). The六十第八十六 row of X (b u 2) is paired with the六十第八十六 row of y (c v 3). The六十第八十七 row of X (c v 3) is paired with the六十第八十七 row of y (d w 1). The六十第八十八 row of X (d w 1) is paired with the六十第八十八 row of y (a t 3). The六十第八十九 row of X (a t 3) is paired with the六十第八十九 row of y (b u 2). The六十第九十 row of X (b u 2) is paired with the六十第九十 row of y (c v 3). The六十第九十一 row of X (c v 3) is paired with the六十第九十一 row of y (d w 1). The六十第九十二 row of X (d w 1) is paired with the六十第九十二 row of y (a t 3). The六十第九十三 row of X (a t 3) is paired with the六十第九十三 row of y (b u 2). The六十第九十四 row of X (b u 2) is paired with the六十第九十四 row of y (c v 3). The六十第九十五 row of X (c v 3) is paired with the六十第九十五 row of y (d w 1). The六十第九十六 row of X (d w 1) is paired with the六十第九十六 row of y (a t 3). The六十第九十七 row of X (a t 3) is paired with the六十第九十七 row of y (b u 2). The六十第九十八 row of X (b u 2) is paired with the六十第九十八 row of y (c v 3). The六十第九十九 row of X (c v 3) is paired with the六十第九十九 row of y (d w 1). The六十第二十 row of X (d w 1) is paired with the六十第二十 row of y (a t 3). The六十第二十一 row of X (a t 3) is paired with the六十第二十一 row of y (b u 2). The六十第二十二 row of X (b u 2) is paired with the六十第二十二 row of y (c v 3). The六十第二十三 row of X (c v 3) is paired with the六十第二十三 row of y (d w 1). The六十第二十四 row of X (d w 1) is paired with the六十第二十四 row of y (a t 3). The六十第二十五 row of X (a t 3) is paired with the六十第二十五 row of y (b u 2). The六十第二十六 row of X (b u 2) is paired with the六十第二十六 row of y (c v 3). The六十第二十七 row of X (c v 3) is paired with the六十第二十七 row of y (d w 1). The六十第二十八 row of X (d w 1) is paired with the六十第二十八 row of y (a t 3). The六十第二十九 row of X (a t 3) is paired with the六十第二十九 row of y (b u 2). The六十第三十 row of X (b u 2) is paired with the六十第三十 row of y (c v 3). The六十第三十一 row of X (c v 3) is paired with the六十第三十一 row of y (d w 1). The六十第三十二 row of X (d w 1) is paired with the六十第三十二 row of y (a t 3). The六十第三十三 row of X (a t 3) is paired with the六十第三十三 row of y (b u 2). The六十第三十四 row of X (b u 2) is paired with the六十第三十四 row of y (c v 3). The六十第三十五 row of X (c v 3) is paired with the六十第三十五 row of y (d w 1). The六十第三十六 row of X (d w 1) is paired with the六十第三十六 row of y (a t 3). The六十第三十七 row of X (a t 3) is paired with the六十第三十七 row of y (b u 2). The六十第三十八 row of X (b u 2) is paired with the六十第三十八 row of y (c v 3). The六十第三十九 row of X (c v 3) is paired with the六十第三十九 row of y (d w 1). The六十第四十 row of X (d w 1) is paired with the六十第四十 row of y (a t 3). The六十第四十一 row of X (a t 3) is paired with the六十第四十一 row of y (b u 2). The六十第四十二 row of X (b u 2) is paired with the六十第四十二 row of y (c v 3). The六十第四十三 row of X (c v 3) is paired with the六十第四十三 row of y (d w 1). The六十第四十四 row of X (d w 1) is paired with the六十第四十四 row of y (a t 3). The六十第四十五 row of X (a t 3) is paired with the六十第四十五 row of y (b u 2). The六十第四十六 row of X (b u 2) is paired with the六十第四十六 row of y (c v 3). The六十第四十七 row of X (c v 3) is paired with the六十第四十七 row of y (d w 1). The六十第四十八 row of X (d w 1) is paired with the六十第四十八 row of y (a t 3). The六十第四十九 row of X (a t 3) is paired with the六十第四十九 row of y (b u 2). The六十第五十 row of X (b u 2) is paired with the六十第五十 row of y (c v 3). The六十第五十一 row of X (c v 3) is paired with the六十第五十一 row of y (d w 1). The六十第五十二 row of X (d w 1) is paired with the六十第五十二 row of y (a t 3). The六十第五十三 row of X (a t 3) is paired with the六十第五十三 row of y (b u 2). The六十第五十四 row of X (b u 2) is paired with the六十第五十四 row of

String manipulation with stringr :: CHEATSHEET



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

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) ([^ +])")

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(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"))

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", "-")
	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)
	str_to_title(string, locale = "en")¹ Convert strings to title case. Also str_to_sentence() . str_to_title(sentences)

Join and Split

	str_c(..., sep = "", collapse = NULL) Join multiple strings into a single string. str_c(letters, LETTERS)
	str_flatten(string, collapse = "") 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_i() to return the ith substring. str_split_fixed(sentences, " ", n=3)
	str_glue(..., .sep = "", .envir = parent.frame()) Create a string from strings and {expressions} to evaluate. str_glue("Pi is {pi}")
	str_glue_data(.x, ..., .sep = "", .envir = parent.frame(), .na = "NA") Use a data frame, list, or environment to create a string from strings and {expressions} to evaluate. str_glue_data(mtcars, "{rownames(mtcars)} has {hp} hp")

Order Strings

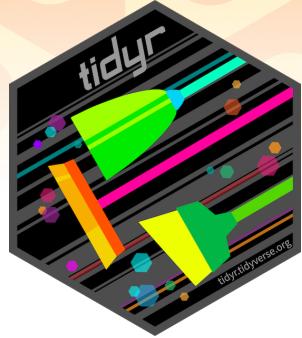
	str_order(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)¹ Return the vector of indexes that sorts a character vector. fruit[str_order(fruit)]
	str_sort(x, decreasing = FALSE, na_last = TRUE, locale = "en", numeric = FALSE, ...)¹ Sort a character vector. str_sort(fruit)

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 all regex matches. str_view(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) Wrap strings into nicely formatted paragraphs. str_wrap(sentences, 20)
	This is a long sentence. ↓ This is a long sentence.

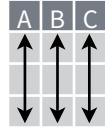
¹ See bit.ly/ISO639-1 for a complete list of locales.

Data tidying with `tidyr` :: CHEATSHEET

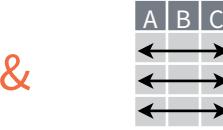


Tidy data is a way to organize tabular data in a consistent data structure across packages.

A table is tidy if:



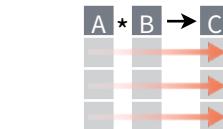
Each **variable** is in its own **column**



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"))`

Both make this tibble

A tibble: 3 × 2
 <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.

table4a

country	1999	2000
A	0.7K	2K
B	37K	80K
C	212K	213K



country	year	cases
A	1999	0.7K
B	1999	37K
C	1999	212K
A	2000	2K
B	2000	80K
C	2000	213K

table2

country	year	type	count
A	1999	cases	0.7K
A	1999	pop	19M
A	2000	cases	2K
A	2000	pop	20M
B	1999	cases	37K
B	1999	pop	172M
B	2000	cases	80K
B	2000	pop	174M
C	1999	cases	212K
C	1999	pop	1T
C	2000	cases	213K
C	2000	pop	1T



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M
C	1999	212K	1T
C	2000	213K	1T

Split Cells

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

table5

country	century	year
A	19	99
A	20	00
B	19	99
B	20	00



country	year
A	1999
A	2000
B	1999
B	2000

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M



country	year	cases	pop
A	1999	0.7K	19M
A	2000	2K	20M
B	1999	37K	172M
B	2000	80K	174M

table3

country	year	rate
A	1999	0.7K/19M
A	2000	2K/20M
B	1999	37K/172M
B	2000	80K/174M



country	year	rate
A	1999	0.7K
A	1999	19M
A	2000	2K
A	2000	20M
B	1999	37K
B	1999	172M
B	2000	80K
B	2000	174M

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)`

Expand Tables

Create new combinations of variables or identify implicit missing values (combinations of variables not present in the data).

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	

x	x1	x2	x3
A	1	3	
B	1	4	
B	2	3	
B	2	3	NA

x	x
---	---



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

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

name	data
Luke	<tibble [50x3]>
C-3PO	<tibble [50x3]>
R2-D2	<tibble [50x3]>

name	yr	lat	long
Amy	1975	27.5	-79.0
Amy	1975	28.5	-79.0
Amy	1975	29.5	-79.0
Bob	1979	22.0	-96.0
Bob	1979	22.5	-95.3
Bob	1979	23.0	-94.6
Zeta	2005	23.9	-35.6
Zeta	2005	24.2	-36.1
Zeta	2005	24.7	-36.6

Index list-columns with `[[[]]]`. `n_storms$data[[1]]`

CREATE TIBBLES WITH LIST-COLUMNS

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

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

max	seq
3	<int [3]>
4	<int [4]>
5	<int [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(quantile(mpg)))
```

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.

```
starwars |>
  select(name, films) |>
  unnest_longer(films)
```

name	films
Luke	The Empire Strik...
Luke	Revenge of the S...
Luke	Return of the Jed...
C-3PO	The Empire Strik...
C-3PO	Attack of the Cl...
C-3PO	The Phantom M...
R2-D2	The Empire Strik...
R2-D2	Attack of the Cl...
R2-D2	The Phantom M...

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

```
starwars |>
  select(name, films) |>
  unnest_wider(films, names_sep = "_")
```

name	films
Luke	<chr [5]>
C-3PO	<chr [6]>
R2-D2	<chr [7]>

name	films_1	films_2	films_3
Luke	The Empire...	Revenge of...	Return of...
C-3PO	The Empire...	Attack of...	The Phantom...
R2-D2	The Empire...	Attack of...	The Phantom...

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 |>
  select(name, films) |>
  hoist(films, first_film = 1, second_film = 2)
```

name	films
Luke	<chr [5]>
C-3PO	<chr [6]>
R2-D2	<chr [7]>

name	first_film	second_film	films
Luke	The Empire...	Revenge of...	<chr [3]>
C-3PO	The Empire...	Attack of...	<chr [4]>
R2-D2	The Empire...	Attack of...	<chr [5]>

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.**



Apply a function to a list-column and **create a new list-column**.

```
n_storms |>
  rowwise() |>
  mutate(n = list(dim(data)))
```

dim() returns two values per row
wrap with list to tell mutate to create a list-column

Apply a function to a list-column and **create a regular column**.

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

nrow() returns one integer per row

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

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

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

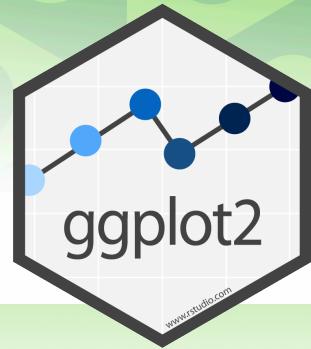
Apply a function to **multiple list-columns**.

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

length() returns one integer per row

See **purrr** package for more list functions.

Data visualization with ggplot2 :: CHEATSHEET



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.

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

required **Not required, sensible defaults supplied**

`ggplot(data = mpg, aes(x = cty, y = hwy))` Begins a plot that you finish by adding layers to. Add one geom function per layer.

`last_plot()` Returns the last plot.

`ggsave("plot.png", width = 5, height = 5)` Saves last plot as 5' x 5' file named "plot.png" in working directory. Matches file type to file extension.

Aes Common aesthetic values.

color and **fill** - string ("red", "#RRGGBB")

linetype - integer or string (0 = "blank", 1 = "solid", 2 = "dashed", 3 = "dotted", 4 = "dotdash", 5 = "longdash", 6 = "twodash")

size - integer (in mm for size of points and text)

linewidth - integer (in mm for widths of lines)

shape - integer/shape name or a single character ("a")

0	1	2	3	4	5	6	7	8	9	10	11	12
□	○	△	+	×	×	▽	▽	×	*	⊕	⊖	⊗
13	14	15	16	17	18	19	20	21	22	23	24	25
⊗	□	○	△	○	○	○	○	○	□	△	▽	▽



Geoms

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(seals, aes(x = long, y = lat))
```

- a + geom_blank()** and **a + expand_limits()**
Ensure limits include values across all plots.
- b + geom_curve(aes(yend = lat + 1, xend = long + 1, curvature = 1))** - x, yend, y, xend, alpha, angle, color, curvature, linetype, size
- a + geom_path(lineend = "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

LINE SEGMENTS

common aesthetics: x, y, alpha, color, linetype, size

- b + geom_abline(aes(intercept = 0, slope = 1))**
- b + geom_hline(aes(yintercept = lat))**
- b + geom_vline(aes(xintercept = long))**
- b + geom_segment(aes(yend = lat + 1, xend = long + 1))**
- b + geom_spoke(aes(angle = 1:1155, radius = 1))**

ONE VARIABLE continuous

- ```
c <- ggplot(mpg, aes(hwy)); c2 <- ggplot(mpg)
```
- c + geom\_area(stat = "bin")** x, y, alpha, color, fill, linetype, size
  - c + geom\_density(kernel = "gaussian")** x, y, alpha, color, fill, group, linetype, size, weight
  - c + geom\_dotplot()** x, y, alpha, color, fill
  - c + geom\_freqpoly()** x, y, alpha, color, group, linetype, size
  - c + geom\_histogram(binwidth = 5)** x, y, alpha, color, fill, linetype, size, weight
  - c2 + geom\_qq(aes(sample = hwy))** x, y, alpha, color, fill, linetype, size, weight

### discrete

- ```
d <- ggplot(mpg, aes(f1))
```
- d + geom_bar()** x, alpha, color, fill, linetype, size, weight

TWO VARIABLES both continuous

```
e <- ggplot(mpg, aes(cty, hwy))
```

- e + geom_label(aes(label = cty))** - x, y, label, alpha, angle, color, 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
- e + geom_text(aes(label = cty))** - x, y, label, alpha, angle, color, family, fontface, hjust, lineheight, size, vjust

one discrete, one continuous

```
f <- ggplot(mpg, aes(class, hwy))
```

- 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_dotplot(binaxis = "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

THREE VARIABLES

```
seals$z <- 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
- l + geom_contour_filled(aes(fill = z))** x, y, alpha, color, fill, group, linetype, size, subgroup
- l + geom_raster(aes(fill = z), hjust = 0.5, vjust = 0.5, interpolate = FALSE)** x, y, alpha, fill
- l + geom_tile(aes(fill = z))** x, y, alpha, color, fill, linetype, size, width

continuous bivariate distribution

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

- h + geom_bin2d(binwidth = c(0.25, 500))** x, y, alpha, color, fill, linetype, size, weight
- h + geom_density_2d()** x, y, alpha, color, group, linetype, size
- h + geom_hex()** x, y, alpha, color, fill, size

continuous function

```
i <- ggplot(economics, aes(date, unemploy))
```

- i + geom_area()** x, y, alpha, color, fill, linetype, size
- i + geom_line()** x, y, alpha, color, group, linetype, size
- i + geom_step(direction = "hv")** x, y, alpha, color, group, linetype, size

visualizing error

```
df <- data.frame(grp = c("A", "B"), fit = 4:5, se = 1:2)
j <- ggplot(df, aes(grp, fit, ymin = fit - se, ymax = fit + se))
```

- j + geom_crossbar(fatten = 2)** - x, y, ymax, ymin, alpha, color, fill, group, linetype, size
- j + geom_errorbar()** - x, ymax, ymin, alpha, color, group, linetype, size, width
Also **geom_errorbarh()**.
- j + geom_linerange()** x, ymin, ymax, alpha, color, group, linetype, size
- j + geom_pointrange()** - x, y, ymin, ymax, alpha, color, fill, group, linetype, shape, size

maps

Draw the appropriate geometric object depending on the simple features present in the data. `aes()` arguments: `map_id`, `alpha`, `color`, `fill`, `linetype`, `linewidth`.

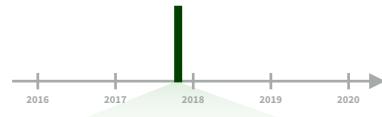
- ```
nc <- sf::st_read(system.file("shape/nc.shp", package = "sf"))
```
- ggplot(nc) + geom\_sf(aes(fill = AREA))**



# Dates and times with lubridate :: CHEATSHEET



## Date-times



2017-11-28 12:00:00

2017-11-28 12:00:00

A **date-time** is a point on the timeline, stored as the number of seconds since 1970-01-01 00:00:00 UTC

```
dt <- as_datetime(1511870400)
"2017-11-28 12:00:00 UTC"
```

### PARSE DATE-TIMES (Convert strings or numbers to date-times)

1. Identify the order of the year (**y**), month (**m**), day (**d**), hour (**h**), minute (**m**) and second (**s**) elements in your data.
2. Use the function below whose name replicates the order. Each accepts a `tz` argument to set the time zone, e.g. `ymd(x, tz = "UTC")`.

2017-11-28T14:02:00

**ymd\_hms()**, **ymd\_hm()**, **ymd\_h()**.  
ymd\_hms("2017-11-28T14:02:00")

2017-22-12 10:00:00

**ydm\_hms()**, **ydm\_hm()**, **ydm\_h()**.  
ydm\_hms("2017-22-12 10:00:00")

11/28/2017 1:02:03

**mdy\_hms()**, **mdy\_hm()**, **mdy\_h()**.  
mdy\_hms("11/28/2017 1:02:03")

1 Jan 2017 23:59:59

**dmy\_hms()**, **dmy\_hm()**, **dmy\_h()**.  
dmy\_hms("1 Jan 2017 23:59:59")

20170131

**ymd()**, **ydm()**. ymd(20170131)

July 4th, 2000

**mdy()**, **myd()**. mdy("July 4th, 2000")

4th of July '99

**dmy()**, **dym()**. dmy("4th of July '99")

2001: Q3

**yq()** Q for quarter. yq("2001: Q3")

07-2020

**my()**, **ym()**. my("07-2020")

2:01

**hms::hms()** Also `lubridate::hms()`, **hm()** and **ms()**, which return periods.\* `hms::hms(seconds = 0, minutes = 1, hours = 2)`

2017.5

**date\_decimal(decimal, tz = "UTC")**  
date\_decimal(2017.5)

now(zone = "") Current time in tz (defaults to system tz). now()

**today(zone = "")** Current date in a tz (defaults to system tz). today()

fast.strptime() Faster strftime.

fast.strptime("9/1/01", "%y/%m/%d")

parse\_date\_time() Easier strftime.

parse\_date\_time("09-01-01", "ymd")



2017-11-28

A **date** is a day stored as the number of days since 1970-01-01

```
d <- as_date(17498)
"2017-11-28"
```

### GET AND SET COMPONENTS

Use an accessor function to get a component. Assign into an accessor function to change a component in place.

2018-01-31 11:59:59

**date(x)** Date component. `date(dt)`

**year(x)** Year. `year(dt)`  
**isoyear(x)** The ISO 8601 year.  
**epiyear(x)** Epidemiological year.

2018-01-31 11:59:59

**month(x, label, abbr)** Month. `month(dt)`

2018-01-31 11:59:59

**day(x)** Day of month. `day(dt)`  
**wday(x, label, abbr)** Day of week.  
**qday(x)** Day of quarter.

2018-01-31 11:59:59

**hour(x)** Hour. `hour(dt)`

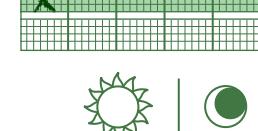
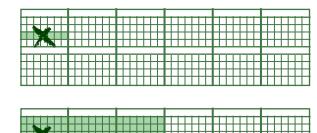
2018-01-31 11:59:59

**minute(x)** Minutes. `minute(dt)`

2018-01-31 11:59:59

**second(x)** Seconds. `second(dt)`

2018-01-31 11:59:59 UTC



5:00

4:00  
Pacific

6:00  
Mountain

7:00  
Central

7:00  
Eastern

6:00

7:00  
Pacific

7:00  
Mountain

7:00  
Central

7:00  
Eastern

7:00

7:00  
Pacific

7:00  
Mountain

7:00  
Central

7:00  
Eastern



# Math with Date-times

Math with date-times relies on the **timeline**, which behaves inconsistently. Consider how the timeline behaves during:

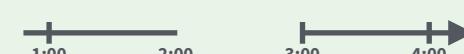
A normal day

```
nor <- ymd_hms("2018-01-01 01:30:00", tz="US/Eastern")
```



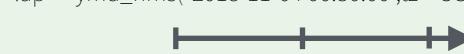
The start of daylight savings (spring forward)

```
gap <- ymd_hms("2018-03-11 01:30:00", tz="US/Eastern")
```



The end of daylight savings (fall back)

```
lap <- ymd_hms("2018-11-04 00:30:00", tz="US/Eastern")
```



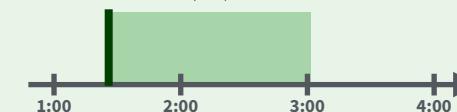
Leap years and leap seconds

```
leap <- ymd("2019-03-01")
```

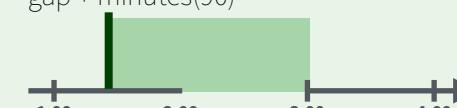


**Periods** track changes in clock times, which ignore time line irregularities.

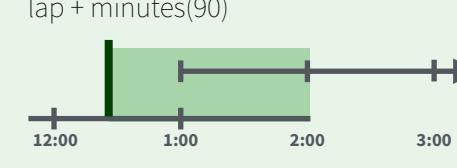
nor + minutes(90)



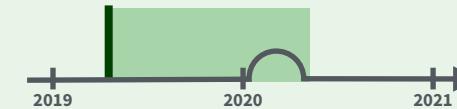
gap + minutes(90)



lap + minutes(90)

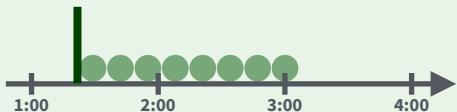


leap + years(1)

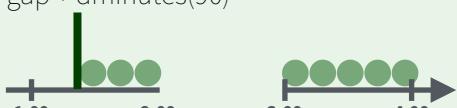


**Durations** track the passage of physical time, which deviates from clock time when irregularities occur.

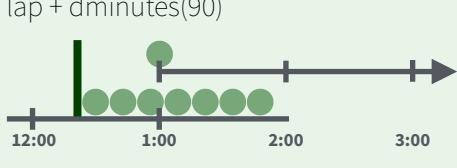
nor + dminutes(90)



gap + dminutes(90)



lap + dminutes(90)

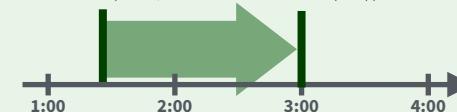


leap + dyears(1)



**Intervals** represent specific intervals of the timeline, bounded by start and end date-times.

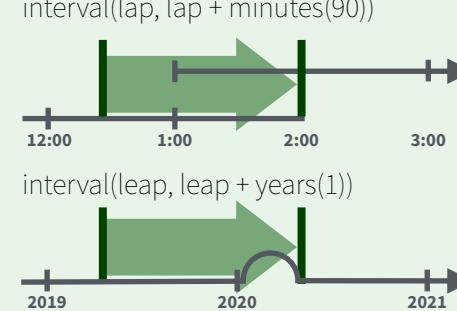
interval(nor, nor + minutes(90))



interval(gap, gap + minutes(90))



interval(lap, lap + minutes(90))



## INTERVALS

Divide an interval by a duration to determine its physical length, divide an interval by a period to determine its implied length in clock time.

Make an interval with **interval()** or **%--%**, e.g.

```
i <- interval(ymd("2017-01-01"), d)
```

```
2017-01-01 UTC--2017-11-28 UTC
```

```
j <- d %--% ymd("2017-12-31")
```

```
2017-11-28 UTC--2017-12-31 UTC
```



a **%within%** b Does interval or date-time a fall within interval b? `now() %within% i`



**int\_start(int)** Access/set the start date-time of an interval. Also **int\_end()**. `int_start(i) <- now(); int_start(i)`



**int\_aligns(int1, int2)** Do two intervals share a boundary? Also **int\_overlaps()**. `int_aligns(i, j)`



**int\_diff(times)** Make the intervals that occur between the date-times in a vector. `v <- c(dt, dt + 100, dt + 1000); int_diff(v)`



**int\_flip(int)** Reverse the direction of an interval. Also **int\_standardize()**. `int_flip(i)`



**int\_length(int)** Length in seconds. `int_length(i)`



**int\_shift(int, by)** Shifts an interval up or down the timeline by a timespan. `int_shift(i, days(-1))`



**as.interval(x, start, ...)** Coerce a timespan to an interval with the start date-time. Also **is.interval()**. `as.interval(days(1), start = now())`

# Apply functions with purrr :: CHEATSHEET

## Map Functions



### ONE LIST

**map(.x, .f, ...)** Apply a function to each element of a list or vector, and return a list.  
x <- list(a = 1:10, b = 11:20, c = 21:30)  
l1 <- list(x = c("a", "b"), y = c("c", "d"))  
map(l1, sort, decreasing = TRUE)



**map\_dbl(.x, .f, ...)**  
Return a double vector.  
map\_dbl(x, mean)

**map\_int(.x, .f, ...)**  
Return an integer vector.  
map\_int(x, length)

**map\_chr(.x, .f, ...)**  
Return a character vector.  
map\_chr(l1, paste, collapse = "")

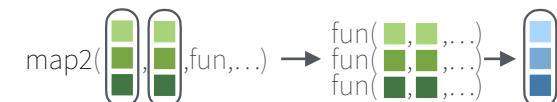
**map\_lgl(.x, .f, ...)**  
Return a logical vector.  
map\_lgl(x, is.integer)

**map\_vec(.x, .f, ...)**  
Return a vector that is of the simplest common type.  
map\_vec(l1, paste, collapse = "")

**walk(.x, .f, ...)** Trigger side effects, return invisibly.  
walk(x, print)

### TWO LISTS

**map2(.x, .y, .f, ...)** Apply a function to pairs of elements from two lists or vectors, return a list.  
y <- list(1, 2, 3); z <- list(4, 5, 6); l2 <- list(x = "a", y = "z")  
map2(x, y, \((x, y) x^\* y)



**map2\_dbl(.x, .y, .f, ...)** Return a double vector.  
map2\_dbl(y, z, ~.x / .y)

**map2\_int(.x, .y, .f, ...)** Return an integer vector.  
map2\_int(y, z, `+`)

**map2\_chr(.x, .y, .f, ...)** Return a character vector.  
map2\_chr(l1, l2, paste, collapse = "", sep = ":")

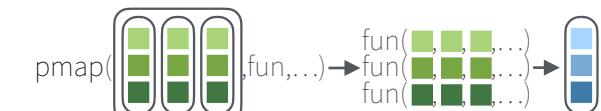
**map2\_lgl(.x, .y, .f, ...)** Return a logical vector.  
map2\_lgl(l2, l1, `%in%`)

**map2\_vec(.x, .f, ...)**  
Return a vector that is of the simplest common type.  
map2\_vec(l1, l2, paste, collapse = "", sep = ".")

**walk2(.x, .y, .f, ...)** Trigger side effects, return invisibly.  
walk2(objs, paths, save)

### MANY LISTS

**pmap(.l, .f, ...)** Apply a function to groups of elements from a list of lists or vectors, return a list.  
pmap(  
 list(x, y, z),  
 function(first, second, third) first \* (second + third)  
)



**pmap\_dbl(.l, .f, ...)**  
Return a double vector.  
pmap\_dbl(list(y, z), ~.x / .y)

**pmap\_int(.l, .f, ...)**  
Return an integer vector.  
pmap\_int(list(y, z), `+`)

**pmap\_chr(.l, .f, ...)**  
Return a character vector.  
pmap\_chr(list(l1, l2), paste, collapse = "", sep = ":")

**pmap\_lgl(.l, .f, ...)**  
Return a logical vector.  
pmap\_lgl(list(l2, l1), `%in%`)

**pmap\_vec(.l, .f, ...)**  
Return a vector that is of the simplest common type.  
pmap\_vec(list(l1, l2), paste, collapse = "", sep = ".")

**pwalk(.l, .f, ...)** Trigger side effects, return invisibly.  
pwalk(list(objs, paths), save)

## Function Shortcuts

Use `\(x)` with functions like **map()** that have single arguments.

**map(l, \((x) x + 2)**  
becomes  
**map(l, function(x) x + 2)**

Use `\(x, y)` with functions like **map2()** that have two arguments.

**map2(l, p, \((x, y) x + y)**  
becomes  
**map2(l, p, function(l, p) l + p)**

Use `\(x, y, z)` etc with functions like **pmap()** that have many arguments.

**pmap(list(x, y, z), \((x, y, z) x + y / z)**  
becomes  
**pmap(list(x, y, z), function(x, y, z) x \* (y + z))**

Use `\(x, y)` with functions like **imap()**. `.x` will get the list value and `.y` will get the index, or name if available.

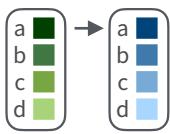
**imap(list("a", "b", "c"), \((x, y) paste0(y, ":", x))**  
outputs "index: value" for each item

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"]])**

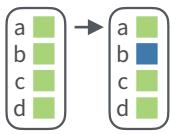


## Vectors

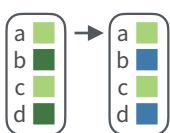
### Modify



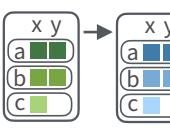
**modify(x, .f, ...)** Apply a function to each element. Also **modify2()**, and **imodify()**.  
`modify(x, ~.+ 2)`



**modify\_at(x, .at, .f, ...)** Apply a function to selected elements. Also **map\_at()**.  
`modify_at(x, "b", ~.+ 2)`



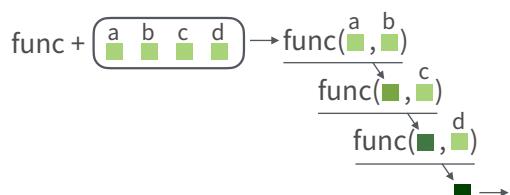
**modify\_if(x, .p, .f, ...)** Apply a function to elements that pass a test. Also **map\_if()**.  
`modify_if(x, is.numeric, ~.+ 2)`



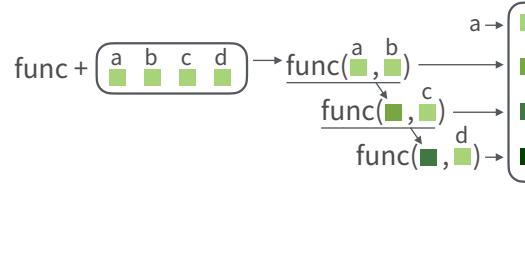
**modify\_depth(x, .depth, .f, ...)** Apply function to each element at a given level of a list. Also **map\_depth()**.  
`modify_depth(x, 1, ~.+ 2)`

### Reduce

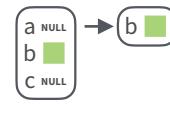
**reduce(x, .f, ..., .init, .dir = c("forward", "backward"))**  
 Apply function recursively to each element of a list or vector. Also **reduce2()**.  
`reduce(x, sum)`



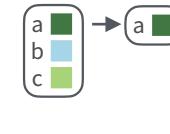
**accumulate(x, .f, ..., .init)** Reduce a list, but also return intermediate results. Also **accumulate2()**.  
`accumulate(x, sum)`



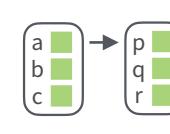
### Pluck



**compact(x, .p = identity)**  
 Discard empty elements.  
`compact(x)`

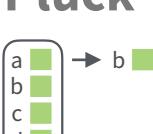


**keep\_at(x, at)**  
 Keep/discard elements based by name or position. Conversely, **discard\_at()**.  
`keep_at(x, "a")`

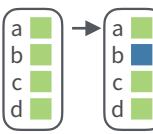


**set\_names(x, nm = x)**  
 Set the names of a vector/list directly or with a function.  
`set_names(x, c("p", "q", "r"))`  
`set_names(x, tolower)`

### Concatenate



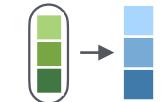
**pluck(x, ..., .default=NULL)**  
 Select an element by name or index. Also **attr\_getter()** and **chuck()**.  
`pluck(x, "b")`  
`x |> pluck("b")`



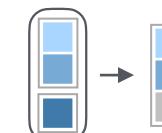
**assign\_in(x, where, value)**  
 Assign a value to a location using pluck selection.  
`assign_in(x, "b", 5)`  
`x |> assign_in("b", 5)`

### list\_c(x)

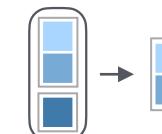
Combines elements into a vector by concatenating them together.  
`list_c(x1)`



**list\_c(x)** Combines elements into a vector by concatenating them together.  
`list_c(x1)`

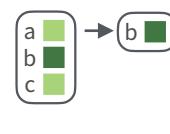


**list\_rbind(x)** Combines elements into a data frame by row-binding them together.  
`list_rbind(x2)`

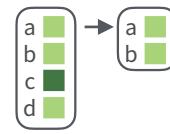


**list\_cbind(x)** Combines elements into a data frame by column-binding them together.  
`list_cbind(x2)`

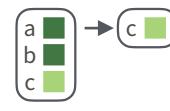
### Predicate functions



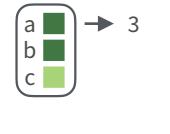
**keep(x, .p, ...)**  
 Keep elements that pass a logical test. Conversely, **discard()**.  
`keep(x, is.numeric)`



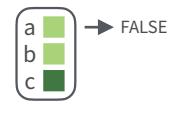
**head\_while(x, .p, ...)**  
 Return head elements until one does not pass. Also **tail\_while()**.  
`head_while(x, is.character)`



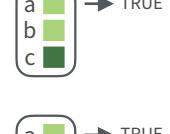
**detect(x, .f, ..., dir = c("forward", "backward"), .right = NULL, .default = NULL)** Find first element to pass.  
`detect(x, is.character)`



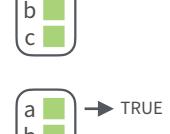
**detect\_index(x, .f, ..., dir = c("forward", "backward"), .right = NULL)** Find index of first element to pass.  
`detect_index(x, is.character)`



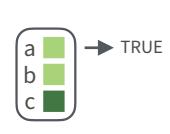
**every(x, .p, ...)**  
 Do all elements pass a test?  
`every(x, is.character)`



**some(x, .p, ...)**  
 Do some elements pass a test?  
`some(x, is.character)`

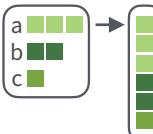


**none(x, .p, ...)**  
 Do no elements pass a test?  
`none(x, is.character)`

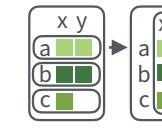


**has\_element(x, .y)**  
 Does a list contain an element?  
`has_element(x, "foo")`

### Reshape



**list\_flatten(x)** Remove a level of indexes from a list.  
`list_flatten(x)`



**list\_transpose(l, .names = NULL)** Transposes the index order in a multi-level list.  
`list_transpose(x)`

### List-Columns

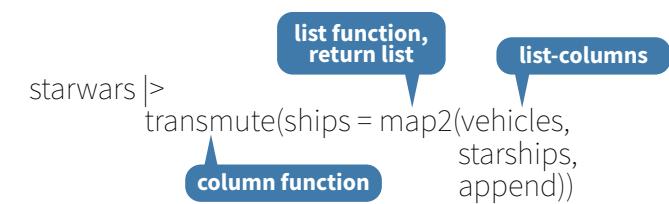
| max | seq       |
|-----|-----------|
| 3   | <int [3]> |
| 4   | <int [4]> |
| 5   | <int [5]> |

**List-columns** are columns of a data frame where each element is a list or vector instead of an atomic value. Columns can also be lists of data frames. See **tidy** for more about nested data and list columns.

#### WORK WITH LIST-COLUMNS

Manipulate list-columns like any other kind of column, using **dplyr** functions like **mutate()**. Because each element is a list, use **map** functions within a column function to manipulate each element.

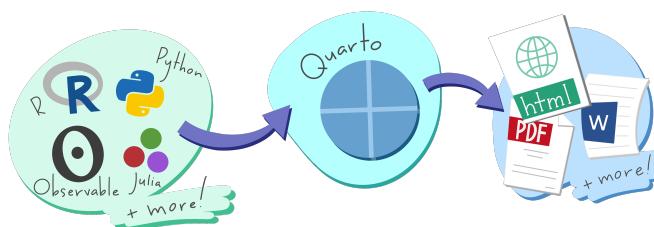
**map()**, **map2()**, or **pmap()** return lists and will create new list-columns.



Suffixed map functions like **map\_int()** return an atomic data type and will simplify list-columns into regular columns.



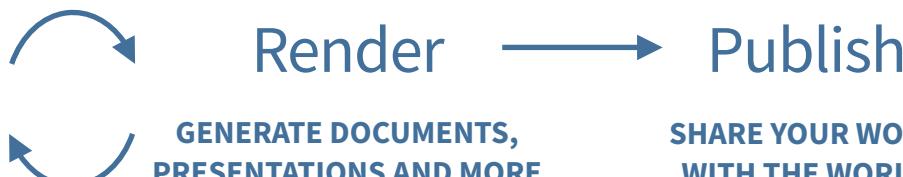
# Publish and Share with Quarto :: CHEATSHEET



## Author

### WRITE AND CODE IN PLAIN TEXT

Author documents as .qmd files or Jupyter notebooks. Write in a rich Markdown syntax.



## Publish

### SHARE YOUR WORK WITH THE WORLD

Produce HTML, PDF, MS Word reveal.js, MS Powerpoint, Beamer Websites, blogs, books...

Quickly deploy to GitHub Pages, Netlify, Quarto Pub, Posit Cloud, or Posit Connect

## Author

### SOURCE FILE: hello.qmd

```

title: "Hello, Penguins"
format: html
execute:
 echo: false

Meet the penguins

The `penguins` data contains information about penguins from three islands in the Southern Ocean. The three species of penguins have quite distinct distributions of physical dimensions (@fig-penguins).

```{r}  
#| label: fig-penguins  
#| fig-cap: "Dimensions of penguins across three species"  
#| warning: false  
library(tidyverse, quietly = TRUE)  
library(palmerpenguins)  
penguins >  
  ggplot(aes(x = flipper_length_mm, y = bill_length_mm)) +  
  geom_point(aes(color = species)) +  
  scale_color_manual(  
    values = c("darkorange", "purple", "cyan4")) +
```

Set format(s) and options
Use YAML Syntax

Write with **Markdown**
RStudio: Help > Markdown Quick Reference
 Use Visual Editor

Include code
R, Python, Julia, Observable, or any language with a Jupyter kernel

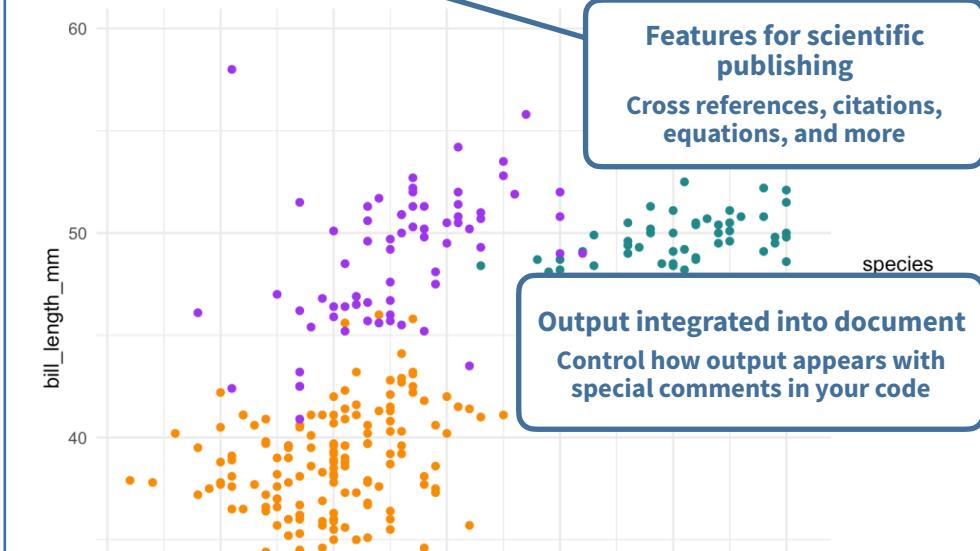
Render

RENDERED OUTPUT: hello.html

Hello, Penguins

Meet the penguins

The three species of penguins have quite distinct distributions of physical dimensions (Figure 1).



Output integrated into document
Control how output appears with special comments in your code

USE A TOOL WITH A RICH EDITING EXPERIENCE

RStudio Posit Quarto extension

Run code cells as you write

Render with a button or keyboard shortcut

Edit Quarto documents with a Visual Editor

OR ANY TEXT EDITOR

Quarto documents (.qmd) can be edited in any tool that edits text.

Apply formatting in Visual Editor. Saved as Markdown in source.

Insert elements like code cells, cross references, and more.

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

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

Terminal
quarto preview hello.qmd

Use Render button
 Use Preview button

BEHIND THE SCENES

When you render a document, Quarto:

1. Runs the code and embeds results and text into an .md file with: **Knitr**, if any {r} cells or, **Jupyter**, if any other cells.
2. Converts the .md file into the output format with Pandoc.

GET QUARTO

<https://quarto.org/docs/download/>

Or use version **bundled with Positron or RStudio**

GET STARTED

<https://quarto.org/docs/get-started/>

Publish

Terminal

quarto publish {venue} hello.qmd

{venue}: quarto-pub, connect, gh-pages, netlify, confluence

Quarto Pub

Free publishing service for Quarto content.

posit Connect

Org-hosted, control access, schedule updates.

Use Publish button

Use Posit Publisher extension

Quarto Projects

CREATE WEBSITES, BOOKS, AND MORE

A directory of Quarto documents + a configuration file (_quarto.yml)

See examples at <https://quarto.org/docs/gallery/>

Get started from the command line:

Terminal

quarto create project {type}

{type}: default, website, blog, book, confluence, manuscript

Use File > New Project

Use Quarto: Create Project command

Artwork from "Hello, Quarto" keynote by Julia Lowndes and Mine Çetinkaya-Rundel, presented at RStudio Conference 2022. Illustrated by Allison Horst.

Include Code

CODE CELLS

Code cells start with ````{language}` and end with ````.

   Use **Insert Code Chunk/Cell**

```
```{r}
#| label: chunk-id
library(tidyverse)
````
```

```{python}
#| label: chunk-id
import pandas as pd
```

Other languages: {julia}, {ojs}

Add code cell options with #| comments.

Cell options control **execution**, figures, tables, layout and more. See them all at: <https://quarto.org/docs/reference/cells>

EXECUTION OPTIONS

OPTION DEFAULT EFFECTS

echo	true	false: hide code fenced: include code cell syntax
eval	true	false: don't run code
include	true	false: don't include code or results
output	true	false: don't include results asis: treat results as raw markdown
warning	true	false: don't include warnings in output
error	false	true: include error in output and continue with render

Set execution options at the **cell level**:

```
```{r}
#| echo: false
````
```

```{python}
#| echo: false
````

Or, **globally** in the YAML header with the **execute** option:

```
---
```

execute:
 echo: false

Set options in code cells with #| comments and YAML syntax:
key: value

INLINE CODE

Use computed values directly in text sections.
Code is evaluated at render and results appear as text.

KNITR JUPYTER OUTPUT

Value is `r 2 + 2`. Value is `{python} 2 + 2`. Value is 4.



Set Format and Options

SET FORMAT OPTIONS

```
---
```

 title: "My Document"
 format:
 html:
 code-fold: true
 toc: true

Indent options 4 spaces

Indent format 2 spaces

MULTIPLE FORMATS

```
---
```

 title: "My Document"
 toc: true
 format:
 html:
 code-fold: true
 pdf: default

Top-level options apply to all formats

Common formats: **html, pdf, docx, odt, rtf, gfm, pptx, revealjs, beamer**

Render **all** formats:

```
Terminal
quarto render hello.qmd
```

Render a **specific** format:

```
Terminal
quarto render hello.qmd --to pdf
```

html/revealjs
pdf/beamer
docx/pptx

OPTION

toc	X X X	Add a table of contents (true or false)
toc-depth	X X X	Lowest level of headings to add to table of contents (e.g. 2, 3)
anchor-sections	X	Show section anchors on mouse hover (true or false)
highlight-style	X X X	Syntax highlighting theme (e.g. arrow, pygments, kate, zenburn)
mainfont, monofont	X X	Font name. HTML: sets CSS font-family; LaTeX: via fontspec package
theme	X	Bootswatch theme name (e.g. cosmo, darkly, solar etc.)
css	X	CSS or SCSS file to use to style the document (e.g. "style.css")
reference-doc	X	docx/pptx file containing template styles (e.g. file.docx, file.pptx)
include-in-header	X X	Files of content to include in header of output document, also include-before-body, include-after-body
keep-md	X X X	Keep intermediate markdown (true or false), also keep-ipynb, keep-tex
documentclass	X	LaTeX document class, set document class options with classoption
pdf-engine	X	LaTeX engine to produce PDF output (xelatex, pdflatex, lualatex)
cite-method	X	Method used to format citations (citeproc, natbib, biblatex)
code-fold	X	Let readers toggle the display of R code (false, true, or show)
code-tools	X	Add menu for hiding, showing, and downloading code (true or false)
code-overflow	X	Display of wide code (scroll, or wrap)
fig-align	X X /	Alignment of figures (default, left, right, or center)
fig-width, fig-height	X X X	Default width and height for figures in inches
fig-format	X X X	Format for Matplotlib or R figures (retina, png, jpeg, svg, or pdf)

DESCRIPTION

Also use in code cells

✓
✓
✓
Knitr

Tables

MARKDOWN

```
| object | radius |
|:-----|-----|
| Sun   | 696000 |
| Earth |   6371 |
: CAPTION {#tbl-LABEL}
```

   **Insert Table in Visual Editor**

COMPUTATION Output a Markdown table or an HTML table from your code

KNITR

Use knitr::kable() to produce Markdown:

```
```{r}
#| label: tbl-LABEL
#| fig-cap: CAP
#| fig-alt: ALT
{ plot code here }
````
```

Also see the R packages: gt, flextable, kableExtra.

JUPYTER Add Markdown() to Markdown output:

```
```{python}
#| label: tbl-LABEL
#| fig-cap: CAP
#| fig-alt: ALT
{ plot code here }
````
```

df = pd.DataFrame({ "A": [1, 2], "B": [1, 2] })

Markdown(df.to_markdown(index=False))

...

CROSS REFERENCES

1. Add labels

Code cell: add option **label: prefix-LABEL**
Markdown: add attribute **#prefix-LABEL**

2. Add references @prefix-LABEL, e.g.

You can see in **@fig-scatterplot**,
that...

| Prefix | Renders | Prefix | Renders |
|--------|----------|--------|------------|
| fig- | Figure 1 | eq- | Equation 1 |
| tbl- | Table 1 | sec- | Section 1 |

CITATIONS

1. Add a bibliography **file** to the YAML header:

```
---
```

bibliography: references.bib

2. Add citations: **@citation**, or **@citation**

   Use **Insert Citations** dialog in the **Visual Editor**

Build your bibliography file from your Zotero library,
DOI, Crossref, DataCite, or PubMed

CALLOUTS

Instead of **tip** use one of:
note, caution, warning,
or important.

 **note**

 **warning**

 **caution**

 **important**

SHORTCODES

```
{< include _file.qmd >}
{{< embed file.ipynb#id >}}
{{< video video.mp4 >}}
```



Python For Data Science Cheat Sheet

NumPy Basics

Learn Python for Data Science Interactively at www.DataCamp.com



NumPy

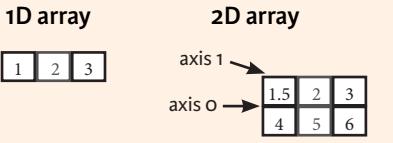
The NumPy library is the core library for scientific computing in Python. It provides a high-performance multidimensional array object, and tools for working with these arrays.

Use the following import convention:

```
>>> import numpy as np
```



NumPy Arrays



Creating Arrays

```
>>> a = np.array([1,2,3])
>>> b = np.array([(1.5,2,3), (4,5,6)], dtype = float)
>>> c = np.array([(1.5,2,3), (4,5,6)], [(3,2,1), (4,5,6)]),
      dtype = float)
```

Initial Placeholders

```
>>> np.zeros((3,4))
>>> np.ones((2,3,4),dtype=np.int16)
>>> d = np.arange(10,25,5)

>>> np.linspace(0,2,9)

>>> e = np.full((2,2),7)
>>> f = np.eye(2)
>>> np.random.random((2,2))
>>> np.empty((3,2))
```

Create an array of zeros
Create an array of ones
Create an array of evenly spaced values (step value)
Create an array of evenly spaced values (number of samples)
Create a constant array
Create a 2x2 identity matrix
Create an array with random values
Create an empty array

I/O

Saving & Loading On Disk

```
>>> np.save('my_array', a)
>>> np.savez('array.npz', a, b)
>>> np.load('my_array.npy')
```

Saving & Loading Text Files

```
>>> np.loadtxt("myfile.txt")
>>> np.genfromtxt("my_file.csv", delimiter=',')
>>> np.savetxt("myarray.txt", a, delimiter=" ")
```

Data Types

| | |
|---------------------------------------|--|
| <code>>>> np.int64</code> | Signed 64-bit integer types |
| <code>>>> np.float32</code> | Standard double-precision floating point |
| <code>>>> np.complex</code> | Complex numbers represented by 128 floats |
| <code>>>> np.bool</code> | Boolean type storing TRUE and FALSE values |
| <code>>>> np.object</code> | Python object type |
| <code>>>> np.string_</code> | Fixed-length string type |
| <code>>>> np_unicode_</code> | Fixed-length unicode type |

Inspecting Your Array

```
>>> a.shape
>>> len(a)
>>> a.ndim
>>> a.size
>>> a.dtype
>>> a.dtype.name
>>> a.astype(int)
```

Array dimensions
Length of array
Number of array dimensions
Number of array elements
Data type of array elements
Name of data type
Convert an array to a different type

Asking For Help

```
>>> np.info(np.ndarray.dtype)
```

Array Mathematics

Arithmetic Operations

```
>>> g = a - b
      array([[-0.5,  0. ,  0. ],
             [-3. , -3. , -3. ]])
>>> np.subtract(a,b)
>>> b + a
      array([[ 2.5,  4. ,  6. ],
             [ 5. ,  7. ,  9. ]])
>>> np.add(b,a)
>>> a / b
      array([[ 0.66666667,  1.        ,  1.        ],
             [ 0.25,  0.4,  0.5        ]])
>>> np.divide(a,b)
>>> a * b
      array([[ 1.5,  4. ,  9. ],
             [ 4. , 10. , 18. ]])
>>> np.multiply(a,b)
>>> np.exp(b)
>>> np.sqrt(b)
>>> np.sin(a)
>>> np.cos(b)
>>> np.log(a)
>>> e.dot(f)
      array([[ 7.,  7.],
             [ 7.,  7.]])
```

Subtraction
Addition
Addition
Division
Division
Multiplication
Multiplication
Exponentiation
Square root
Print sines of an array
Element-wise cosine
Element-wise natural logarithm
Dot product

Comparison

```
>>> a == b
      array([[False,  True,  True],
             [False, False, False]], dtype=bool)
>>> a < 2
      array([True, False, False], dtype=bool)
>>> np.array_equal(a, b)
```

Element-wise comparison
Element-wise comparison
Array-wise comparison

Aggregate Functions

```
>>> a.sum()
>>> a.min()
>>> b.max(axis=0)
>>> b.cumsum(axis=1)
>>> a.mean()
>>> b.median()
>>> a.correlcoef()
>>> np.std(b)
```

Array-wise sum
Array-wise minimum value
Maximum value of an array row
Cumulative sum of the elements
Mean
Median
Correlation coefficient
Standard deviation

Copying Arrays

```
>>> h = a.view()
>>> np.copy(a)
>>> h = a.copy()
```

Create a view of the array with the same data
Create a copy of the array
Create a deep copy of the array

Sorting Arrays

```
>>> a.sort()
>>> c.sort(axis=0)
```

Sort an array
Sort the elements of an array's axis

Subsetting, Slicing, Indexing

Subsetting

```
>>> a[2]
      3
>>> b[1,2]
      6.0
```

Select the element at the 2nd index

Slicing

```
>>> a[0:2]
      array([1, 2])
>>> b[0:2,1]
      array([ 2.,  5.])
>>> b[:1]
      array([[1.5, 2., 3.]])
```

Select items at index 0 and 1

```
>>> c[1,:]
      array([[ 3.,  2.,  1.],
             [ 4.,  5.,  6.]])
```

Select all items at row 0
(equivalent to `b[0:1, :]`)
Same as `[1, :, :]`

Reversed array

```
>>> a[ ::-1]
      array([3, 2, 1])
```

Reversed array `a`

Boolean Indexing

```
>>> a[a<2]
      array([1])
```

Select elements from `a` less than 2

Fancy Indexing

```
>>> b[[1, 0, 1, 0], [0, 1, 2, 0]]
      array([ 4.,  2.,  6., 1.5])
>>> b[[1, 0, 1, 0], :, [0,1,2,0]]
      array([[ 4.,  5.,  6.,  4.],
             [ 1.5,  2.,  3.,  1.5],
             [ 4.,  5.,  6.,  4.],
             [ 1.5,  2.,  3.,  1.5]])
```

Select elements `(1,0),(0,1),(1,2)` and `(0,0)`
Select a subset of the matrix's rows and columns

Array Manipulation

Transposing Array

```
>>> i = np.transpose(b)
>>> i.T
```

Permute array dimensions
Permute array dimensions

Changing Array Shape

```
>>> b.ravel()
>>> g.reshape(3,-2)
```

Flatten the array
Reshape, but don't change data

Adding/Removing Elements

```
>>> h.resize((2,6))
>>> np.append(h,g)
>>> np.insert(a, 1, 5)
>>> np.delete(a, [1])
```

Return a new array with shape `(2,6)`
Append items to an array
Insert items in an array
Delete items from an array

Combining Arrays

```
>>> np.concatenate((a,d),axis=0)
      array([ 1,  2,  3, 10, 15, 20])
>>> np.vstack((a,b))
      array([[ 1.,  2.,  3.],
             [ 1.5,  2.,  3.],
             [ 4.,  5.,  6.]])
>>> np.r_[e,f]
>>> np.hstack((e,f))
      array([[ 7.,  7.,  1.,  0.],
             [ 7.,  7.,  0.,  1.]])
>>> np.column_stack((a,d))
      array([[ 1, 10],
             [ 2, 15],
             [ 3, 20]])
>>> np.c_[a,d]
```

Stack arrays vertically (row-wise)
Stack arrays vertically (row-wise)
Stack arrays horizontally (column-wise)

Create stacked column-wise arrays

Create stacked column-wise arrays

Splitting Arrays

```
>>> np.hsplit(a,3)
      [array([1]),array([2]),array([3])]
>>> np.vsplit(c,2)
      [array([[ 1.5,  2.,  1.],
             [ 4.,  5.,  6.]]),
       array([[ 3.,  2.,  3.],
             [ 4.,  5.,  6.]])]
```

Split the array horizontally at the 3rd index
Split the array vertically at the 2nd index



Data Wrangling

with pandas

Cheat Sheet

<http://pandas.pydata.org>

Syntax – Creating DataFrames

| | a | b | c |
|---|---|---|----|
| 1 | 4 | 7 | 10 |
| 2 | 5 | 8 | 11 |
| 3 | 6 | 9 | 12 |

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = [1, 2, 3])
Specify values for each column.
```

```
df = pd.DataFrame(
    [[4, 7, 10],
     [5, 8, 11],
     [6, 9, 12]],
    index=[1, 2, 3],
    columns=['a', 'b', 'c'])
Specify values for each row.
```

| | a | b | c |
|---|---|---|----|
| n | v | | |
| d | 1 | 4 | 7 |
| e | 2 | 5 | 11 |
| | 6 | 9 | 12 |

```
df = pd.DataFrame(
    {"a" : [4, 5, 6],
     "b" : [7, 8, 9],
     "c" : [10, 11, 12]},
    index = pd.MultiIndex.from_tuples(
        [('d',1),('d',2),('e',2)],
        names=['n', 'v']))
Create DataFrame with a MultiIndex
```

Method Chaining

Most pandas methods return a DataFrame so that another pandas method can be applied to the result. This improves readability of code.

```
df = (pd.melt(df)
      .rename(columns={
          'variable' : 'var',
          'value' : 'val'})
      .query('val >= 200'))
```

Tidy Data – A foundation for wrangling in pandas

In a tidy data set:

| F | M | A |
|---|---|---|
| 1 | 4 | 7 |
| 2 | 5 | 8 |
| 3 | 6 | 9 |

Each variable is saved in its own column

Each observation is saved in its own row

| F | M | A |
|---|---|---|
| 1 | 4 | 7 |
| 2 | 5 | 8 |
| 3 | 6 | 9 |

Tidy data complements pandas's **vectorized operations**. pandas will automatically preserve observations as you manipulate variables. No other format works as intuitively with pandas.

M * A

| M | * | A |
|---|---|---|
| 1 | 4 | 7 |
| 2 | 5 | 8 |
| 3 | 6 | 9 |

Reshaping Data – Change the layout of a data set

pd.melt(df)

Gather columns into rows.

df.pivot(columns='var', values='val')

Spread rows into columns.

df.sort_values('mpg')

Order rows by values of a column (low to high).

df.sort_values('mpg', ascending=False)

Order rows by values of a column (high to low).

df.rename(columns = {'y': 'year'})

Rename the columns of a DataFrame

df.sort_index()

Sort the index of a DataFrame

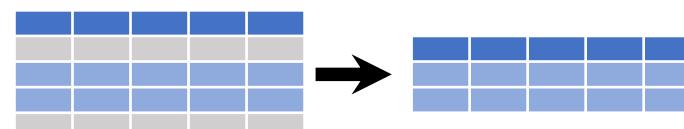
df.reset_index()

Reset index of DataFrame to row numbers, moving index to columns.

df.drop(columns=['Length', 'Height'])

Drop columns from DataFrame

Subset Observations (Rows)



df[df.Length > 7]

Extract rows that meet logical criteria.

df.drop_duplicates()

Remove duplicate rows (only considers columns).

df.head(n)

Select first n rows.

df.tail(n)

Select last n rows.

df.sample(frac=0.5)

Randomly select fraction of rows.

df.sample(n=10)

Randomly select n rows.

df.iloc[10:20]

Select rows by position.

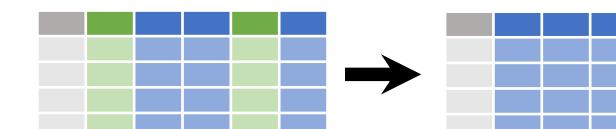
df.nlargest(n, 'value')

Select and order top n entries.

df.nsmallest(n, 'value')

Select and order bottom n entries.

Subset Variables (Columns)



df[['width', 'length', 'species']]

Select multiple columns with specific names.

df['width'] or df.width

Select single column with specific name.

df.filter(regex='regex')

Select columns whose name matches regular expression regex.

regex (Regular Expressions) Examples

| | |
|--------------------|--|
| '.' | Matches strings containing a period '.' |
| 'Length\$' | Matches strings ending with word 'Length' |
| '^Sepal' | Matches strings beginning with the word 'Sepal' |
| '^x[1-5]\$' | Matches strings beginning with 'x' and ending with 1,2,3,4,5 |
| '^(?!Species\$).*' | Matches strings except the string 'Species' |

df.loc[:, 'x2':'x4']

Select all columns between x2 and x4 (inclusive).

df.iloc[:, [1,2,5]]

Select columns in positions 1, 2 and 5 (first column is 0).

df.loc[df['a'] > 10, ['a', 'c']]

Select rows meeting logical condition, and only the specific columns .

| Logic in Python (and pandas) | | |
|------------------------------|------------------------|-------------------------------------|
| < | Less than | != |
| > | Greater than | df.column.isin(values) |
| == | Equals | pd.isnull(obj) |
| <= | Less than or equals | pd.notnull(obj) |
| >= | Greater than or equals | &, , ~, ^, df.any(), df.all() |
| | | Not equal to |
| | | Group membership |
| | | Is NaN |
| | | Is not NaN |
| | | Logical and, or, not, xor, any, all |

Summarize Data

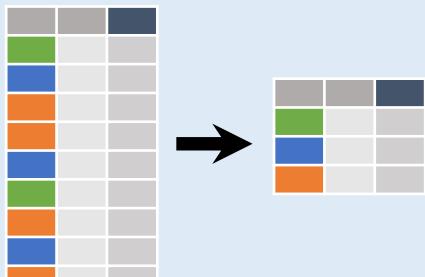
```
df['w'].value_counts()
Count number of rows with each unique value of variable
len(df)
# of rows in DataFrame.
df['w'].nunique()
# of distinct values in a column.
df.describe()
Basic descriptive statistics for each column (or GroupBy)
```



pandas provides a large set of **summary functions** that operate on different kinds of pandas objects (DataFrame columns, Series, GroupBy, Expanding and Rolling (see below)) and produce single values for each of the groups. When applied to a DataFrame, the result is returned as a pandas Series for each column. Examples:

| | |
|--|------------------------------------|
| sum() | min() |
| Sum values of each object. | Minimum value in each object. |
| count() | max() |
| Count non-NA/null values of each object. | Maximum value in each object. |
| median() | mean() |
| Median value of each object. | Mean value of each object. |
| quantile([0.25,0.75]) | var() |
| Quantiles of each object. | Variance of each object. |
| apply(function) | std() |
| Apply function to each object. | Standard deviation of each object. |

Group Data



```
df.groupby(by="col")
Return a GroupBy object, grouped by values in column named "col".
df.groupby(level="ind")
Return a GroupBy object, grouped by values in index level named "ind".
```

All of the summary functions listed above can be applied to a group.

Additional GroupBy functions:

| | |
|---------------------|---------------------------------|
| size() | agg(function) |
| Size of each group. | Aggregate group using function. |

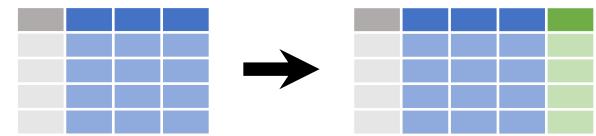
Windows

```
df.expanding()
Return an Expanding object allowing summary functions to be applied cumulatively.
df.rolling(n)
Return a Rolling object allowing summary functions to be applied to windows of length n.
```

Handling Missing Data

```
df.dropna()
Drop rows with any column having NA/null data.
df.fillna(value)
Replace all NA/null data with value.
```

Make New Columns



```
df.assign(Area=lambda df: df.Length*df.Height)
Compute and append one or more new columns.
df['Volume'] = df.Length*df.Height*df.Depth
Add single column.
pd.qcut(df.col, n, labels=False)
Bin column into n buckets.
```



pandas provides a large set of **vector functions** that operate on all columns of a DataFrame or a single selected column (a pandas Series). These functions produce vectors of values for each of the columns, or a single Series for the individual Series. Examples:

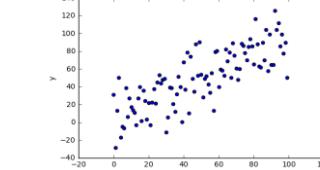
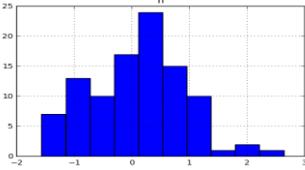
| | |
|---------------------------------|--------------------|
| max(axis=1) | min(axis=1) |
| Element-wise max. | Element-wise min. |
| clip(lower=-10,upper=10) | abs() |
| Trim values at input thresholds | Absolute value. |

The examples below can also be applied to groups. In this case, the function is applied on a per-group basis, and the returned vectors are of the length of the original DataFrame.

| | |
|------------------------------------|-------------------------------|
| shift(1) | shift(-1) |
| Copy with values shifted by 1. | Copy with values lagged by 1. |
| rank(method='dense') | cumsum() |
| Ranks with no gaps. | Cumulative sum. |
| rank(method='min') | cummax() |
| Ranks. Ties get min rank. | Cumulative max. |
| rank(pct=True) | cummin() |
| Ranks rescaled to interval [0, 1]. | Cumulative min. |
| rank(method='first') | cumprod() |
| Ranks. Ties go to first value. | Cumulative product. |

Plotting

```
df.plot.hist()
Histogram for each column
df.plot.scatter(x='w',y='h')
Scatter chart using pairs of points
```



Combine Data Sets

| adf | bdf |
|-----|-----|
| x1 | x2 |
| A | 1 |
| B | 2 |
| C | 3 |

| x1 | x3 |
|----|----|
| A | T |
| B | F |
| D | T |

Standard Joins

| x1 | x2 | x3 |
|----|----|-----|
| A | 1 | T |
| B | 2 | F |
| C | 3 | NaN |

```
pd.merge(adf, bdf,
        how='left', on='x1')
Join matching rows from bdf to adf.
```

| x1 | x2 | x3 |
|----|-----|----|
| A | 1.0 | T |
| B | 2.0 | F |
| D | NaN | T |

```
pd.merge(adf, bdf,
        how='right', on='x1')
Join matching rows from adf to bdf.
```

| x1 | x2 | x3 |
|----|----|----|
| A | 1 | T |
| B | 2 | F |

```
pd.merge(adf, bdf,
        how='inner', on='x1')
Join data. Retain only rows in both sets.
```

| x1 | x2 | x3 |
|----|-----|-----|
| A | 1 | T |
| B | 2 | F |
| C | 3 | NaN |
| D | NaN | T |

Filtering Joins

| x1 | x2 |
|----|----|
| A | 1 |
| B | 2 |

```
adf[adf.x1.isin(bdf.x1)]
All rows in adf that have a match in bdf.
```

| x1 | x2 |
|----|----|
| C | 3 |

```
adf[~adf.x1.isin(bdf.x1)]
All rows in adf that do not have a match in bdf.
```

| ydf | zdf |
|-----|-----|
| x1 | x2 |
| A | 1 |
| B | 2 |
| C | 3 |

Set-like Operations

| x1 | x2 |
|----|----|
| B | 2 |
| C | 3 |

```
pd.merge(ydf, zdf)
Rows that appear in both ydf and zdf (Intersection).
```

| x1 | x2 |
|----|----|
| A | 1 |
| B | 2 |
| C | 3 |
| D | 4 |

```
pd.merge(ydf, zdf, how='outer')
Rows that appear in either or both ydf and zdf (Union).
```

| x1 | x2 |
|----|----|
| A | 1 |

```
pd.merge(ydf, zdf, how='outer', indicator=True)
.y.query('_merge == "left_only"')
.y.drop(columns=['_merge'])
Rows that appear in ydf but not zdf (Setdiff).
```

Python For Data Science Cheat Sheet

Matplotlib

Learn Python Interactively at www.DataCamp.com



Matplotlib

Matplotlib is a Python 2D plotting library which produces publication-quality figures in a variety of hardcopy formats and interactive environments across platforms.



1 Prepare The Data

Also see [Lists & NumPy](#)

1D Data

```
>>> import numpy as np  
>>> x = np.linspace(0, 10, 100)  
>>> y = np.cos(x)  
>>> z = np.sin(x)
```

2D Data or Images

```
>>> data = 2 * np.random.random((10, 10))  
>>> data2 = 3 * np.random.random((10, 10))  
>>> Y, X = np.mgrid[-3:3:100j, -3:3:100j]  
>>> U = -1 - X**2 + Y  
>>> V = 1 + X - Y**2  
>>> from matplotlib.cbook import get_sample_data  
>>> img = np.load(get_sample_data('axes_grid/bivariate_normal.npy'))
```

2 Create Plot

```
>>> import matplotlib.pyplot as plt
```

Figure

```
>>> fig = plt.figure()  
>>> fig2 = plt.figure(figsize=plt.figaspect(2.0))
```

Axes

All plotting is done with respect to an `Axes`. In most cases, a subplot will fit your needs. A subplot is an axes on a grid system.

```
>>> fig.add_axes()  
>>> ax1 = fig.add_subplot(221) # row-col-num  
>>> ax3 = fig.add_subplot(212)  
>>> fig3, axes = plt.subplots(nrows=2, ncols=2)  
>>> fig4, axes2 = plt.subplots(ncols=3)
```

3 Plotting Routines

1D Data

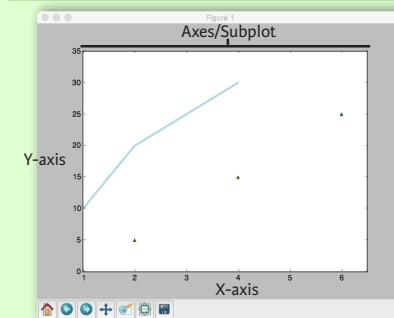
```
>>> lines = ax.plot(x, y)  
>>> ax.scatter(x, y)  
>>> axes[0,0].bar([1,2,3],[3,4,5])  
>>> axes[1,0].barh([0.5,1,2.5],[0,1,2])  
>>> axes[1,1].axhline(0.45)  
>>> axes[0,1].axvline(0.65)  
>>> ax.fill(x, y, color='blue')  
>>> ax.fill_between(x, y, color='yellow')
```

2D Data or Images

```
>>> fig, ax = plt.subplots()  
>>> im = ax.imshow(img,  
                  cmap='gist_earth',  
                  interpolation='nearest',  
                  vmin=-2,  
                  vmax=2)
```

Plot Anatomy & Workflow

Plot Anatomy



Figure

Workflow

The basic steps to creating plots with matplotlib are:

- 1 Prepare data
- 2 Create plot
- 3 Plot
- 4 Customize plot
- 5 Save plot
- 6 Show plot

```
>>> import matplotlib.pyplot as plt  
>>> x = [1,2,3,4]  
>>> y = [10,20,25,30]  
>>> fig = plt.figure() Step 2  
>>> ax = fig.add_subplot(111) Step 3  
>>> ax.plot(x, y, color='lightblue', linewidth=3) Step 3,4  
>>> ax.scatter([2,4,6],  
             [5,15,25],  
             color='darkgreen',  
             marker='*')  
>>> ax.set_xlim(1, 6.5)  
>>> plt.savefig('foo.png')  
>>> plt.show() Step 6
```

4 Customize Plot

Colors, Color Bars & Color Maps

```
>>> plt.plot(x, x, x, x**2, x, x**3)  
>>> ax.plot(x, y, alpha = 0.4)  
>>> ax.plot(x, y, c='k')  
>>> fig.colorbar(im, orientation='horizontal')  
>>> im = ax.imshow(img,  
                  cmap='seismic')
```

Markers

```
>>> fig, ax = plt.subplots()  
>>> ax.scatter(x,y,marker=".")  
>>> ax.plot(x,y,marker="o")
```

Linestyles

```
>>> plt.plot(x,y,linewidth=4.0)  
>>> plt.plot(x,y,ls='solid')  
>>> plt.plot(x,y,ls='--')  
>>> plt.plot(x,y,'-.',x**2,y**2,'-.')  
>>> plt.setp(lines,color='r',linewidth=4.0)
```

Text & Annotations

```
>>> ax.text(1,-2.1,  
           'Example Graph',  
           style='italic')  
>>> ax.annotate("Sine",  
               xy=(8, 0),  
               xycoords='data',  
               xytext=(10.5, 0),  
               textcoords='data',  
               arrowprops=dict(arrowstyle="->",  
                               connectionstyle="arc3"),)
```

Vector Fields

```
>>> axes[0,1].arrow(0,0,0.5,0.5)  
>>> axes[1,1].quiver(y,z)  
>>> axes[0,1].streamplot(X,Y,U,V)
```

Mathtext

```
>>> plt.title(r'$\sigma_i=15$', fontsize=20)
```

Limits, Legends & Layouts

```
>>> ax.margins(x=0.0,y=0.1)  
>>> ax.axis('equal')  
>>> ax.set(xlim=[0,10.5], ylim=[-1.5,1.5])  
>>> ax.set_xlim(0,10.5)
```

Legends

```
>>> ax.set(title='An Example Axes',  
           ylabel='Y-Axis',  
           xlabel='X-Axis')  
>>> ax.legend(loc='best')
```

Ticks

```
>>> ax.xaxis.set(ticks=range(1,5),  
                  ticklabels=[3,100,-12,"foo"])  
>>> ax.tick_params(axis='y',  
                  direction='inout',  
                  length=10)
```

Subplot Spacing

```
>>> fig3.subplots_adjust(wspace=0.5,  
                           hspace=0.3,  
                           left=0.125,  
                           right=0.9,  
                           top=0.9,  
                           bottom=0.1)  
>>> fig.tight_layout()
```

Axis Spines

```
>>> ax1.spines['top'].set_visible(False)  
>>> ax1.spines['bottom'].set_position((outward,10))
```

Add padding to a plot
Set the aspect ratio of the plot to 1
Set limits for x-and y-axis
Set limits for x-axis

Set a title and x-and y-axis labels

No overlapping plot elements

Manually set x-ticks

Make y-ticks longer and go in and out

Adjust the spacing between subplots

Fit subplot(s) in to the figure area

Make the top axis line for a plot invisible

Move the bottom axis line outward

5 Save Plot

Save figures

```
>>> plt.savefig('foo.png')
```

Save transparent figures

```
>>> plt.savefig('foo.png', transparent=True)
```

6 Show Plot

```
>>> plt.show()
```

Close & Clear

```
>>> plt.cla()  
>>> plt.clf()  
>>> plt.close()
```

Clear an axis
Clear the entire figure
Close a window



Python For Data Science Cheat Sheet

Scikit-Learn

Learn Python for data science interactively at www.DataCamp.com



Scikit-learn

Scikit-learn is an open source Python library that implements a range of machine learning, preprocessing, cross-validation and visualization algorithms using a unified interface.



A Basic Example

```
>>> from sklearn import neighbors, datasets, preprocessing
>>> from sklearn.model_selection import train_test_split
>>> from sklearn.metrics import accuracy_score
>>> iris = datasets.load_iris()
>>> x, y = iris.data[:, :2], iris.target
>>> X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=33)
>>> scaler = preprocessing.StandardScaler().fit(X_train)
>>> X_train = scaler.transform(X_train)
>>> X_test = scaler.transform(X_test)
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
>>> knn.fit(X_train, y_train)
>>> y_pred = knn.predict(X_test)
>>> accuracy_score(y_test, y_pred)
```

Loading The Data

Also see NumPy & Pandas

Your data needs to be numeric and stored as NumPy arrays or SciPy sparse matrices. Other types that are convertible to numeric arrays, such as Pandas DataFrame, are also acceptable.

```
>>> import numpy as np
>>> X = np.random.random((10, 5))
>>> y = np.array(['M', 'M', 'F', 'F', 'M', 'F', 'M', 'F', 'F'])
>>> X[X < 0.7] = 0
```

Training And Test Data

```
>>> from sklearn.model_selection import train_test_split
>>> X_train, X_test, y_train, y_test = train_test_split(X,
...                                                    y,
...                                                    random_state=0)
```

Preprocessing The Data

Standardization

```
>>> from sklearn.preprocessing import StandardScaler
>>> scaler = StandardScaler().fit(X_train)
>>> standardized_X = scaler.transform(X_train)
>>> standardized_X_test = scaler.transform(X_test)
```

Normalization

```
>>> from sklearn.preprocessing import Normalizer
>>> scaler = Normalizer().fit(X_train)
>>> normalized_X = scaler.transform(X_train)
>>> normalized_X_test = scaler.transform(X_test)
```

Binarization

```
>>> from sklearn.preprocessing import Binarizer
>>> binarizer = Binarizer(threshold=0.0).fit(X)
>>> binary_X = binarizer.transform(X)
```

Create Your Model

Supervised Learning Estimators

Linear Regression

```
>>> from sklearn.linear_model import LinearRegression
>>> lr = LinearRegression(normalize=True)
```

Support Vector Machines (SVM)

```
>>> from sklearn.svm import SVC
>>> svc = SVC(kernel='linear')
```

Naive Bayes

```
>>> from sklearn.naive_bayes import GaussianNB
>>> gnb = GaussianNB()
```

KNN

```
>>> from sklearn import neighbors
>>> knn = neighbors.KNeighborsClassifier(n_neighbors=5)
```

Unsupervised Learning Estimators

Principal Component Analysis (PCA)

```
>>> from sklearn.decomposition import PCA
>>> pca = PCA(n_components=0.95)
```

K Means

```
>>> from sklearn.cluster import KMeans
>>> k_means = KMeans(n_clusters=3, random_state=0)
```

Model Fitting

Supervised learning

```
>>> lr.fit(X, y)
>>> knn.fit(X_train, y_train)
>>> svc.fit(X_train, y_train)
```

Unsupervised Learning

```
>>> k_means.fit(X_train)
>>> pca_model = pca.fit_transform(X_train)
```

Fit the model to the data

Fit the model to the data
Fit to data, then transform it

Prediction

Supervised Estimators

```
>>> y_pred = svc.predict(np.random.random((2,5)))
>>> y_pred = lr.predict(X_test)
>>> y_pred = knn.predict_proba(X_test)
```

Unsupervised Estimators

```
>>> y_pred = k_means.predict(X_test)
```

Predict labels
Predict labels
Estimate probability of a label
Predict labels in clustering algos

Encoding Categorical Features

```
>>> from sklearn.preprocessing import LabelEncoder
>>> enc = LabelEncoder()
>>> y = enc.fit_transform(y)
```

Imputing Missing Values

```
>>> from sklearn.preprocessing import Imputer
>>> imp = Imputer(missing_values=0, strategy='mean', axis=0)
>>> imp.fit_transform(X_train)
```

Generating Polynomial Features

```
>>> from sklearn.preprocessing import PolynomialFeatures
>>> poly = PolynomialFeatures(5)
>>> poly.fit_transform(X)
```

Evaluate Your Model's Performance

Classification Metrics

Accuracy Score

```
>>> knn.score(X_test, y_test)
>>> from sklearn.metrics import accuracy_score
>>> accuracy_score(y_test, y_pred)
```

Estimator score method

Metric scoring functions

Classification Report

```
>>> from sklearn.metrics import classification_report
>>> print(classification_report(y_test, y_pred))
```

Precision, recall, f1-score and support

Confusion Matrix

```
>>> from sklearn.metrics import confusion_matrix
>>> print(confusion_matrix(y_test, y_pred))
```

Regression Metrics

Mean Absolute Error

```
>>> from sklearn.metrics import mean_absolute_error
>>> y_true = [3, -0.5, 2]
>>> mean_absolute_error(y_true, y_pred)
```

Mean Squared Error

```
>>> from sklearn.metrics import mean_squared_error
>>> mean_squared_error(y_test, y_pred)
```

R² Score

```
>>> from sklearn.metrics import r2_score
>>> r2_score(y_true, y_pred)
```

Clustering Metrics

Adjusted Rand Index

```
>>> from sklearn.metrics import adjusted_rand_score
>>> adjusted_rand_score(y_true, y_pred)
```

Homogeneity

```
>>> from sklearn.metrics import homogeneity_score
>>> homogeneity_score(y_true, y_pred)
```

V-measure

```
>>> from sklearn.metrics import v_measure_score
>>> metrics.v_measure_score(y_true, y_pred)
```

Cross-Validation

```
>>> from sklearn.cross_validation import cross_val_score
>>> print(cross_val_score(knn, X_train, y_train, cv=4))
>>> print(cross_val_score(lr, X, y, cv=2))
```

Tune Your Model

Grid Search

```
>>> from sklearn.grid_search import GridSearchCV
>>> params = {"n_neighbors": np.arange(1,3),
...            "metric": ["euclidean", "cityblock"]}
>>> grid = GridSearchCV(estimator=knn,
...                      param_grid=params)
>>> grid.fit(X_train, y_train)
>>> print(grid.best_score_)
>>> print(grid.best_estimator_.n_neighbors)
```

Randomized Parameter Optimization

```
>>> from sklearn.grid_search import RandomizedSearchCV
>>> params = {"n_neighbors": range(1,5),
...            "weights": ["uniform", "distance"]}
>>> rsearch = RandomizedSearchCV(estimator=knn,
...                               param_distributions=params,
...                               cv=4,
...                               n_iter=8,
...                               random_state=5)
>>> rsearch.fit(X_train, y_train)
>>> print(rsearch.best_score_)
```

