

Data Transformation With dplyr :: CHEAT SHEET



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



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

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

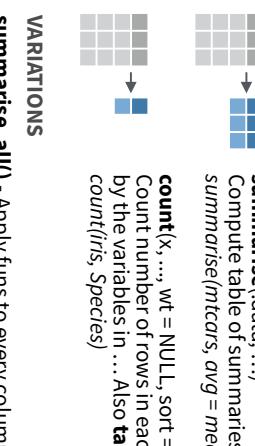
Summarise Cases

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

summary function

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

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



VARIATIONS

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

Group Cases

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



ADD CASES

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

ungroup(**x**, ...)

Returns ungroupped copy of table.
ungroup(g_iris)

Manipulate Cases

EXTRACT CASES

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

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

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

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

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

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

select function

summarise(mtcars, avg = mean(mpg))

count(iris, Species)

group_by(iris, Species)

Logical and boolean operators to use with filter()

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

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

ARRANGE CASES

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

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

ADD CASES

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

ungroup(**x**, ...)

Returns ungroupped copy of table.
ungroup(g_iris)

Manipulate Variables

EXTRACT VARIABLES

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

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

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

use_helpers with **select ()**
e.g. **select(iris, starts_with("Sepal"))**

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

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

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

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

add_column(**tbl**, .cols, .new, .before = NULL, .after = NULL) Add new column(s). Also **add_count()**, **add_column(mtcars, new = 1:32)**, **add_row(faithful, eruptions = 1, waiting = 1)**

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

Vector Functions

Summary Functions

TO USE WITH MUTATE()

mutate() and **transmute()** apply vectorized functions to columns to create new columns.

Vectorized functions take vectors as input and return vectors of the same length as output.

vectorized function



OFFSETS

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

CUMULATIVE AGGREGATES

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

RANKINGS

dplyr::cume_dist() - Proportion of all values <= dplyr::dense_rank() - rank 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

<=,>,>=, != == - logical comparisons

dplyr::between() - x >= left & x <= right

dplyr::near() - safe == for floating point numbers

MISC

dplyr::case_when() - multi-case if_else()
Species == "Versicolor" ~ "Versi",
Species == "virginica" ~ "virgi",
TRUE ~ Species))

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.

dplyr::coalesce() - first non-NA values by element across a set of vectors

dplyr::if_else() - element-wise if + else()
dplyr::is_if() - replace specific values with NA

pmax() - element-wise max()
pmin() - element-wise min()

dplyr::recode() - Vectorized switch()
dplyr::recode_factor() - Vectorized switch() for factors

TO USE WITH SUMMARISE()

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

summary function



COUNTS

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

LOCATION

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

LOGICALS

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

POSITION/ORDER

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

RANK

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

SPREAD

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

Combine Tables

COMBINE CASES

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 3 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 2 |
| b u 2 | b u 3 |
| c v 3 | c v 4 |

COMBINE VARIABLES

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

bind_cols(...)

Returns tables placed side by side as a single table.
BE SURE THAT ROWS ALIGN.

bind_rows(...)

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

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

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

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 2 |
| b u 2 | b u 3 |
| c v 3 | c v 4 |

left_join(x, y, by = NULL, ...)

Join matching values from y to x.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

right_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join matching values from x to y.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

inner_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join data. Retain only rows with matches.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

full_join(x, y, by = NULL, copy = FALSE, suffix=c("x", "y"), ...)

Join data. Retain all values, all rows.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

setdiff(x, y, ...)

Rows that appear in x but not y.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

union(x, y, ...)

Rows that appear in x or y.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

setdiff(x, y, by = NULL, ...)

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

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

semi_join(x, y, by = NULL, ...)

Useful to see what will be joined.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

anti_join(x, y, by = NULL, ...)

Return rows of x that do not have a match in y. Useful to see what will NOT BE JOINED.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

has_rownames()

Match in y. Useful to see what will NOT BE JOINED.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |

remove_rownames()

Not be joined.

| X | Y |
|-----------|-----------|
| A B C | A B C |
| a t 1 | a t 1 |
| b u 2 | b u 2 |
| c v 3 | c v 3 |