

Data transformation with dplyr : : CHEATSHEET



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



&



pipes

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

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

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



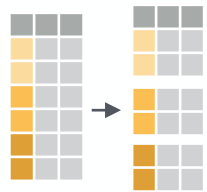
summarize(.data, ...)
Compute table of summaries.
mtcars > summarize(avg = mean(mpg))



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 tidyR cheat sheet for list-column workflow.



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

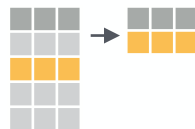
ungroup(x, ...) Returns ungrouped copy of table.

g_mtcars <- mtcars > group_by(cyl)
ungroup(g_mtcars)

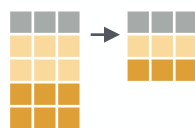
Manipulate Cases

EXTRACT CASES

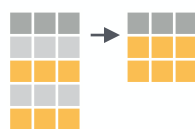
Row functions return a subset of rows as a new table.



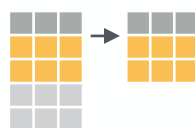
filter(.data, ..., .preserve = FALSE) Extract rows that meet logical criteria.
mtcars > filter(mpg > 20)



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() | ! | & | |

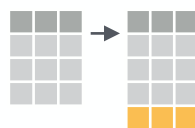
See **?base::Logic** and **?Comparison** for help.

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 CASES



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) | num_range(prefix, range) | ; e.g., mpg:cyl |
| ends_with(match) | all_of(x)/any_of(x, ..., vars) | ! e.g., !gear |
| starts_with(match) | matches(match) | everything() |

MANIPULATE MULTIPLE VARIABLES AT ONCE

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))

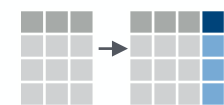


c_across(.cols) Compute across columns in row-wise data.
df >
rowwise() >
mutate(x_total = sum(c_across(1:2)))

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)