Data Transformation with dplyr:: **cheat sheet**

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





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



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).





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



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

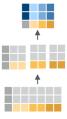
VARIATIONS

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

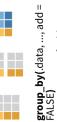
Group Cases

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

mtcars %>%



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



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

Manipulate Cases

EXTRACT CASES

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



weight = NULL, .env = parent.frame()) Randomly select fraction of rows. distinct(.data, ..., keep_all = FALSE) Remove
rows with duplicate values.
distinct(iris, Species) sample_frac(tbl, size = 1, replace = FALSE **^**

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

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)*

Logical and boolean operators to use with filter()

See ?base::logic and ?Comparison for help. %in% is.na() IJ

ARRANGE CASES



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

ADD CASES



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

Manipulate Variables

dplyr

EXTRACT VARIABLES

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



pull(data, var = -1) Extract column values as a vector. Choose by name or index. pull(iris, Sepal.Length) select(.data,...)
Extract columns as a table. Also select_if().
select(iris, Sepal.Length, Species)

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

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

MAKE NEW VARIABLES

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





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

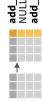
xor()



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



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



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



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



grouped by ... g_iris <- group_by(iris, Species)

Returns copy of table