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



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 <=</pre>
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",
```

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 dplvr::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.



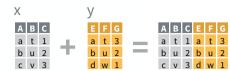


AB tibble::column_to_rownames() 1 a t b t 1 a Move col into row names. 3 c v v 3 c a | > column_to_rownames(var = "C")

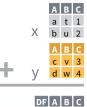
Also tibble::has_rownames() and tibble::remove_rownames().

Combine Tables

COMBINE VARIABLES



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



COMBINE CASES

bind_rows(..., .id = NULL) Returns tables one on top of the other as a single table. Set .id to a column name to add a column of the original table names (as pictured).

RELATIONAL DATA

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



ABCD left_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE,c v 3 NA na_matches = "na") Join matching values from v to x.



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



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



full_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ..., keep = FALSE, c v 3 NA na_matches = "na") Join data. Retain all dw NA 1 values, all rows.

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

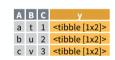


ABC semi_join(x, y, by = NULL, copy = FALSE, ..., na_matches = "na") Return rows of x that have a match in y. Use to see what will be included in a join.



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

Use a "Nest Join" to inner join one table to another into a nested data frame.



nest_join(x, y, by = NULL, copy = FALSE, keep = FALSE, name = NULL, ...) Join data, nesting matches from y in a single new data frame column.

COLUMN MATCHING FOR JOINS



Use by = c("col1", "col2", ...) to specify one or more common columns to match on. left join(x, y, by = "A")



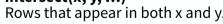
Use a named vector, by = c("col1" = "col2"), to match on columns that have different names in each table. $left_{join}(x, y, by = c("C" = "D"))$

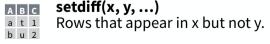


Use **suffix** to specify the suffix to give to unmatched columns that have the same name in both tables. left join(x, y, by = c("C" = "D"),suffix = c("1", "2")

SET OPERATIONS

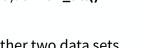








union(x, y, ...) Rows that appear in x or y, duplicates removed). union_all() retains duplicates.



Use **setequal()** to test whether two data sets contain the exact same rows (in any order).

