# Package: dplyr (via r-universe)

June 26, 2024

Type Package

Title A Grammar of Data Manipulation

Version 1.1.4.9000

**Description** A fast, consistent tool for working with data frame like objects, both in memory and out of memory.

License MIT + file LICENSE

URL https://dplyr.tidyverse.org, https://github.com/tidyverse/dplyr

BugReports https://github.com/tidyverse/dplyr/issues

**Depends** R (>= 3.6.0)

**Imports** cli (>= 3.6.2), generics, glue (>= 1.3.2), lifecycle (>= 1.0.3), magrittr (>= 1.5), methods, pillar (>= 1.9.0), R6, rlang (>= 1.1.3), tibble (>= 3.2.0), tidyselect (>= 1.2.0), utils, vctrs (>= 0.6.4)

Suggests bench, broom, callr, covr, DBI, dbplyr (>= 2.2.1), ggplot2, knitr, Lahman, lobstr, microbenchmark, nycflights13, purrr, rmarkdown, RMySQL, RPostgreSQL, RSQLite, stringi (>= 1.7.6), testthat (>= 3.1.5), tidyr (>= 1.3.0), withr

VignetteBuilder knitr

Config/Needs/website tidyverse, shiny, pkgdown, tidyverse/tidytemplate

Config/testthat/edition 3

**Encoding** UTF-8

LazyData true

**Roxygen** list(markdown = TRUE)

RoxygenNote 7.3.1

**Repository** https://tidyverse.r-universe.dev

RemoteUrl https://github.com/tidyverse/dplyr

RemoteRef HEAD

RemoteSha 0005f6768fa765b3bba5148711967d58b6013037

# Contents

across	3
all_vars	7
arrange	8
auto_copy	10
band_members	10
between	11
bind cols	12
bind rows	13
case match	14
case when	16
coalesce	19
compute	20
consecutive id	21
context	22
copy to	23
count	24
cross join	26
cumall	27
	28
desc	20
distinct	30
dnlvr hv	31
evolain	36
filter	30
filter joins	30
alimnse	41
group by	12
group_oy	
group_cois	45
group_map	43
ident	47
ident	40
$\lim_{t \to \infty}  h_{t} ^{2} = \int_{t}  h_{t} ^{2} dt = \int_{t}  h_{t} ^{2} d$	49 50
	50
	54
	55
	38
na_11	00
near	66
nest_join	67
nth	69
nule	71
n_distinct	72
order_by	73
percent_rank	74
pick	75
pull	76

#### across

recode															•									77
reframe														•										80
relocate														•										82
rename	•	•	•		•			•		•			•	•	•			•			•		•	83
rows	•	•	•		•			•		•			•	•	•			•			•		•	85
rowwise		•	•					•						•	•			•		•		•		88
row_number		•	•					•						•	•			•						90
scoped		•	•					•						•					•					91
select		•	•					•						•	•			•						93
setops		•	•					•						•					•					97
slice	•	•	•		•			•		•			•	•	•			•			•		•	98
sql		•	•					•						•					•					102
starwars	•	•	•		•			•		•			•	•	•			•			•		•	102
storms		•	•					•											•					103
summarise	•	•	•		•			•		•			•	•	•			•			•		•	104
tbl		•	•					•											•					107
vars	•		•					•						•	•	•		•	•	•		•		108

# Index

across

Apply a function (or functions) across multiple columns

#### Description

across() makes it easy to apply the same transformation to multiple columns, allowing you to use select() semantics inside in "data-masking" functions like summarise() and mutate(). See vignette("colwise") for more details.

if\_any() and if\_all() apply the same predicate function to a selection of columns and combine the results into a single logical vector: if\_any() is TRUE when the predicate is TRUE for *any* of the selected columns, if\_all() is TRUE when the predicate is TRUE for *all* selected columns.

If you just need to select columns without applying a transformation to each of them, then you probably want to use pick() instead.

across() supersedes the family of "scoped variants" like summarise\_at(), summarise\_if(), and summarise\_all().

# Usage

```
across(.cols, .fns, ..., .names = NULL, .unpack = FALSE)
if_any(.cols, .fns, ..., .names = NULL)
if_all(.cols, .fns, ..., .names = NULL)
```

3

#### Arguments

.cols

.fns

<tidy-select> Columns to transform. You can't select grouping columns because they are already automatically handled by the verb (i.e. summarise() or mutate()).

Functions to apply to each of the selected columns. Possible values are:

- A function, e.g. mean.
- A purrr-style lambda, e.g. ~ mean(.x, na.rm = TRUE)
- A named list of functions or lambdas, e.g. list(mean = mean, n\_miss = ~ sum(is.na(.x)). Each function is applied to each column, and the output is named by combining the function name and the column name using the glue specification in .names.

Within these functions you can use cur\_column() and cur\_group() to access the current column and grouping keys respectively.

# ... [Deprecated]

Additional arguments for the function calls in .fns are no longer accepted in ... because it's not clear when they should be evaluated: once per across() or once per group? Instead supply additional arguments directly in .fns by using a lambda. For example, instead of across(a:b, mean, na.rm = TRUE) write across(a:b, ~ mean(.x, na.rm = TRUE)).

.names A glue specification that describes how to name the output columns. This can use {.col} to stand for the selected column name, and {.fn} to stand for the name of the function being applied. The default (NULL) is equivalent to "{.col}" for the single function case and "{.col}\_{.fn}" for the case where a list is used for .fns.

. unpack [Experimental]

Optionally unpack data frames returned by functions in .fns, which expands the df-columns out into individual columns, retaining the number of rows in the data frame.

- If FALSE, the default, no unpacking is done.
- If TRUE, unpacking is done with a default glue specification of "{outer}\_{inner}".
- Otherwise, a single glue specification can be supplied to describe how to name the unpacked columns. This can use {outer} to refer to the name originally generated by .names, and {inner} to refer to the names of the data frame you are unpacking.

#### Value

across() typically returns a tibble with one column for each column in .cols and each function in .fns. If .unpack is used, more columns may be returned depending on how the results of .fns are unpacked.

if\_any() and if\_all() return a logical vector.

#### Timing of evaluation

R code in dplyr verbs is generally evaluated once per group. Inside across() however, code is evaluated once for each combination of columns and groups. If the evaluation timing is important,

#### across

for example if you're generating random variables, think about when it should happen and place your code in consequence.

```
gdf <-
  tibble(g = c(1, 1, 2, 3), v1 = 10:13, v2 = 20:23) %>%
 group_by(g)
set.seed(1)
# Outside: 1 normal variate
n <- rnorm(1)
gdf %>% mutate(across(v1:v2, ~ .x + n))
#> # A tibble: 4 x 3
#> # Groups: g [3]
#>
        g
             v1 v2
#>
     <dbl> <dbl> <dbl>
#> 1
        1 9.37 19.4
#> 2
        1 10.4
                 20.4
#> 3
        2 11.4
                 21.4
#> 4
        3 12.4
                 22.4
# Inside a verb: 3 normal variates (ngroup)
gdf %>% mutate(n = rnorm(1), across(v1:v2, ~ .x + n))
#> # A tibble: 4 x 4
#> # Groups: g [3]
#>
             v1
                   v2
        g
                           n
     <dbl> <dbl> <dbl> <dbl>
#>
#> 1
        1 10.2 20.2 0.184
#> 2
        1 11.2 21.2 0.184
#> 3
        2 11.2 21.2 -0.836
        3 14.6 24.6 1.60
#> 4
# Inside `across()`: 6 normal variates (ncol * ngroup)
gdf %>% mutate(across(v1:v2, ~ .x + rnorm(1)))
#> # A tibble: 4 x 3
#> # Groups: g [3]
#>
             v1
                   v2
        g
#>
     <dbl> <dbl> <dbl>
#> 1
        1 10.3 20.7
#> 2
        1 11.3 21.7
#> 3
        2 11.2 22.6
        3 13.5 22.7
#> 4
```

## See Also

c\_across() for a function that returns a vector

#### across

# Examples

```
# For better printing
iris <- as_tibble(iris)</pre>
# across() ------
# Different ways to select the same set of columns
# See <https://tidyselect.r-lib.org/articles/syntax.html> for details
iris %>%
 mutate(across(c(Sepal.Length, Sepal.Width), round))
iris %>%
 mutate(across(c(1, 2), round))
iris %>%
 mutate(across(1:Sepal.Width, round))
iris %>%
 mutate(across(where(is.double) & !c(Petal.Length, Petal.Width), round))
# Using an external vector of names
cols <- c("Sepal.Length", "Petal.Width")</pre>
iris %>%
 mutate(across(all_of(cols), round))
# If the external vector is named, the output columns will be named according
# to those names
names(cols) <- tolower(cols)</pre>
iris %>%
 mutate(across(all_of(cols), round))
# A purrr-style formula
iris %>%
 group_by(Species) %>%
 summarise(across(starts_with("Sepal"), ~ mean(.x, na.rm = TRUE)))
# A named list of functions
iris %>%
 group_by(Species) %>%
 summarise(across(starts_with("Sepal"), list(mean = mean, sd = sd)))
# Use the .names argument to control the output names
iris %>%
 group_by(Species) %>%
 summarise(across(starts_with("Sepal"), mean, .names = "mean_{.col}"))
iris %>%
 group_by(Species) %>%
 summarise(across(starts_with("Sepal"), list(mean = mean, sd = sd), .names = "{.col}.{.fn}"))
# If a named external vector is used for column selection, .names will use
# those names when constructing the output names
iris %>%
 group_by(Species) %>%
 summarise(across(all_of(cols), mean, .names = "mean_{.col}"))
# When the list is not named, .fn is replaced by the function's position
```

all\_vars

```
iris %>%
 group_by(Species) %>%
 summarise(across(starts_with("Sepal"), list(mean, sd), .names = "{.col}.fn{.fn}"))
# When the functions in .fns return a data frame, you typically get a
# "packed" data frame back
quantile_df <- function(x, probs = c(0.25, 0.5, 0.75)) {</pre>
 tibble(quantile = probs, value = quantile(x, probs))
}
iris %>%
 reframe(across(starts_with("Sepal"), quantile_df))
# Use .unpack to automatically expand these packed data frames into their
# individual columns
iris %>%
 reframe(across(starts_with("Sepal"), quantile_df, .unpack = TRUE))
# .unpack can utilize a glue specification if you don't like the defaults
iris %>%
 reframe(across(starts_with("Sepal"), quantile_df, .unpack = "{outer}.{inner}"))
# This is also useful inside mutate(), for example, with a multi-lag helper
multilag <- function(x, lags = 1:3) {</pre>
 names(lags) <- as.character(lags)</pre>
 purrr::map_dfr(lags, lag, x = x)
}
iris %>%
 group_by(Species) %>%
 mutate(across(starts_with("Sepal"), multilag, .unpack = TRUE)) %>%
 select(Species, starts_with("Sepal"))
# if_any() and if_all() -------
iris %>%
 filter(if_any(ends_with("Width"), ~ . > 4))
iris %>%
 filter(if_all(ends_with("Width"), ~ . > 2))
```

all\_vars

Apply predicate to all variables

#### Description

#### [Superseded]

all\_vars() and any\_vars() were only needed for the scoped verbs, which have been superseded by the use of across() in an existing verb. See vignette("colwise") for details.

These quoting functions signal to scoped filtering verbs (e.g. filter\_if() or filter\_all()) that a predicate expression should be applied to all relevant variables. The all\_vars() variant takes the

intersection of the predicate expressions with & while the any\_vars() variant takes the union with |.

# Usage

```
all_vars(expr)
any_vars(expr)
```

#### Arguments

expr

<data-masking> An expression that returns a logical vector, using . to refer to
the "current" variable.

## See Also

vars() for other quoting functions that you can use with scoped verbs.

Order rows using column values

#### Description

arrange() orders the rows of a data frame by the values of selected columns.

Unlike other dplyr verbs, arrange() largely ignores grouping; you need to explicitly mention grouping variables (or use .by\_group = TRUE) in order to group by them, and functions of variables are evaluated once per data frame, not once per group.

## Usage

```
arrange(.data, ..., .by_group = FALSE)
## S3 method for class 'data.frame'
arrange(.data, ..., .by_group = FALSE, .locale = NULL)
```

## Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<pre><data-masking> Variables, or functions of variables. Use desc() to sort a variable in descending order.</data-masking></pre>
.by_group	If TRUE, will sort first by grouping variable. Applies to grouped data frames only.
.locale	The locale to sort character vectors in.
	• If NULL, the default, uses the "C" locale unless the dplyr.legacy_locale global option escape hatch is active. See the dplyr-locale help page for more details.

- If a single string from stringi::stri\_locale\_list() is supplied, then this will be used as the locale to sort with. For example, "en" will sort with the American English locale. This requires the stringi package.
- If "C" is supplied, then character vectors will always be sorted in the C locale. This does not require stringi and is often much faster than supplying a locale identifier.

The C locale is not the same as English locales, such as "en", particularly when it comes to data containing a mix of upper and lower case letters. This is explained in more detail on the locale help page under the Default locale section.

## Details

## **Missing values:**

Unlike base sorting with sort(), NA are:

- always sorted to the end for local data, even when wrapped with desc().
- treated differently for remote data, depending on the backend.

# Value

An object of the same type as .data. The output has the following properties:

- All rows appear in the output, but (usually) in a different place.
- · Columns are not modified.
- Groups are not modified.
- Data frame attributes are preserved.

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

## See Also

Other single table verbs: filter(), mutate(), reframe(), rename(), select(), slice(), summarise()

```
arrange(mtcars, cyl, disp)
arrange(mtcars, desc(disp))

# grouped arrange ignores groups
by_cyl <- mtcars %>% group_by(cyl)
by_cyl %>% arrange(desc(wt))
# Unless you specifically ask:
by_cyl %>% arrange(desc(wt), .by_group = TRUE)
```

```
# use embracing when wrapping in a function;
# see ?rlang::args_data_masking for more details
tidy_eval_arrange <- function(.data, var) {
  .data %>%
    arrange({{ var }})
}
tidy_eval_arrange(mtcars, mpg)
# Use `across()` or `pick()` to select columns with tidy-select
iris %>% arrange(pick(starts_with("Sepal")))
iris %>% arrange(across(starts_with("Sepal"), desc))
```

auto\_copy Copy tables to same source, if necessary

# Description

Copy tables to same source, if necessary

#### Usage

auto\_copy(x, y, copy = FALSE, ...)

## Arguments

х, у	y will be copied to x, if necessary.
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.
•••	Other arguments passed on to methods.

|--|

# Description

These data sets describe band members of the Beatles and Rolling Stones. They are toy data sets that can be displayed in their entirety on a slide (e.g. to demonstrate a join).

# Usage

band\_members

band\_instruments

band\_instruments2

# between

## Format

Each is a tibble with two variables and three observations

# Details

band\_instruments and band\_instruments2 contain the same data but use different column names for the first column of the data set. band\_instruments uses name, which matches the name of the key column of band\_members; band\_instruments2 uses artist, which does not.

## Examples

band\_members
band\_instruments
band\_instruments2

```
between
```

Detect where values fall in a specified range

## Description

This is a shortcut for  $x \ge \text{left } \& x \le \text{right}$ , implemented for local vectors and translated to the appropriate SQL for remote tables.

## Usage

between(x, left, right)

# Arguments

xA vectorleft, rightBoundary values. Both left and right are recycled to the size of x.

## Details

x, left, and right are all cast to their common type before the comparison is made.

#### Value

A logical vector the same size as x.

## See Also

join\_by() if you are looking for documentation for the between() overlap join helper.

## Examples

```
between(1:12, 7, 9)
x <- rnorm(1e2)
x[between(x, -1, 1)]
# On a tibble using `filter()`
filter(starwars, between(height, 100, 150))</pre>
```

bind\_cols

Bind multiple data frames by column

# Description

Bind any number of data frames by column, making a wider result. This is similar to do.call(cbind, dfs).

Where possible prefer using a join to combine multiple data frames. bind\_cols() binds the rows in order in which they appear so it is easy to create meaningless results without realising it.

## Usage

```
bind_cols(
    ...,
    .name_repair = c("unique", "universal", "check_unique", "minimal")
)
```

# Arguments

	Data frames to combine. Each argument can either be a data frame, a list that
	could be a data frame, or a list of data frames. Inputs are recycled to the same
	length, then matched by position.
.name_repair	One of "unique", "universal", or "check_unique". See vctrs::vec_as_names()
	for the meaning of these options.

## Value

A data frame the same type as the first element of ....

## Examples

df1 <- tibble(x = 1:3)
df2 <- tibble(y = 3:1)
bind\_cols(df1, df2)</pre>

# Row sizes must be compatible when column-binding
try(bind\_cols(tibble(x = 1:3), tibble(y = 1:2)))

bind\_rows

# Description

Bind any number of data frames by row, making a longer result. This is similar to do.call(rbind, dfs), but the output will contain all columns that appear in any of the inputs.

#### Usage

bind\_rows(..., .id = NULL)

## Arguments

	Data frames to combine. Each argument can either be a data frame, a list that could be a data frame, or a list of data frames. Columns are matched by name, and any missing columns will be filled with NA.
.id	The name of an optional identifier column. Provide a string to create an out- put column that identifies each input. The column will use names if available, otherwise it will use positions.

## Value

A data frame the same type as the first element of ....

```
df1 <- tibble(x = 1:2, y = letters[1:2])
df2 <- tibble(x = 4:5, z = 1:2)
# You can supply individual data frames as arguments:
bind_rows(df1, df2)
# Or a list of data frames:
bind_rows(list(df1, df2))
# When you supply a column name with the `.id` argument, a new
# column is created to link each row to its original data frame
bind_rows(list(df1, df2), .id = "id")
bind_rows(list(a = df1, b = df2), .id = "id")</pre>
```

case\_match

#### Description

This function allows you to vectorise multiple switch() statements. Each case is evaluated sequentially and the first match for each element determines the corresponding value in the output vector. If no cases match, the .default is used.

case\_match() is an R equivalent of the SQL "simple" CASE WHEN statement.

## Connection to case\_when():

While case\_when() uses logical expressions on the left-hand side of the formula, case\_match() uses values to match against .x with. The following two statements are roughly equivalent:

```
case_when(
    x %in% c("a", "b") ~ 1,
    x %in% "c" ~ 2,
    x %in% c("d", "e") ~ 3
)
case_match(
    x,
    c("a", "b") ~ 1,
    "c" ~ 2,
    c("d", "e") ~ 3
)
```

#### Usage

case\_match(.x, ..., .default = NULL, .ptype = NULL)

## Arguments

. X	A vector to match against.
	<dynamic-dots> A sequence of two-sided formulas: old_values ~ new_value. The right hand side (RHS) determines the output value for all values of .x that match the left hand side (LHS).</dynamic-dots>
	The LHS must evaluate to the same type of vector as .x. It can be any length, allowing you to map multiple .x values to the same RHS value. If a value is repeated in the LHS, i.e. a value in .x matches to multiple cases, the first match is used.
	The RHS inputs will be coerced to their common type. Each RHS input will be recycled to the size of .x.
default	The value used when values in .x aren't matched by any of the LHS inputs. If NULL, the default, a missing value will be useddefault is recycled to the size of .x.
.ptype	An optional prototype declaring the desired output type. If not supplied, the output type will be taken from the common type of all RHS inputs and .default.

## case\_match

# Value

A vector with the same size as .x and the same type as the common type of the RHS inputs and .default (if not overridden by .ptype).

## See Also

case\_when()

```
x <- c("a", "b", "a", "d", "b", NA, "c", "e")
# `case_match()` acts like a vectorized `switch()`.
# Unmatched values "fall through" as a missing value.
case_match(
  х,
  ″a″~1,
  "b" ~ 2,
  "c" ~ 3,
  "d" ~ 4
)
# Missing values can be matched exactly, and `.default` can be used to
# control the value used for unmatched values of x^{-1}
case_match(
  х,
  "a" ~ 1,
  "b" ~ 2,
  "c" ~ 3,
  "d" ~ 4,
  NA ~ 0,
  .default = 100
)
# Input values can be grouped into the same expression to map them to the
# same output value
case_match(
  х,
  c("a", "b") ~ "low",
  c("c", "d", "e") ~ "high"
)
# `case_match()` isn't limited to character input:
y <- c(1, 2, 1, 3, 1, NA, 2, 4)
case_match(
  у,
  c(1, 3) ~ "odd",
  c(2, 4) ~ "even",
  .default = "missing"
)
```

```
# Setting `.default` to the original vector is a useful way to replace
# selected values, leaving everything else as is
case_match(y, NA ~ 0, .default = y)
starwars %>%
 mutate(
   # Replace missings, but leave everything else alone
   hair_color = case_match(hair_color, NA ~ "unknown", .default = hair_color),
   # Replace some, but not all, of the species
    species = case_match(
      species,
      "Human" ~ "Humanoid",
      "Droid" ~ "Robot",
     c("Wookiee", "Ewok") ~ "Hairy",
      .default = species
   ),
    .keep = "used"
 )
```

case\_when

A general vectorised if-else

# Description

This function allows you to vectorise multiple if\_else() statements. Each case is evaluated sequentially and the first match for each element determines the corresponding value in the output vector. If no cases match, the .default is used as a final "else" statement.

case\_when() is an R equivalent of the SQL "searched" CASE WHEN statement.

## Usage

case\_when(..., .default = NULL, .ptype = NULL, .size = NULL)

# Arguments

... 
 <dynamic-dots> A sequence of two-sided formulas. The left hand side (LHS) determines which values match this case. The right hand side (RHS) provides the replacement value.
 The LHS inputs must evaluate to logical vectors.
 The RHS inputs will be coerced to their common type.
 All inputs will be recycled to their common size. That said, we encourage all LHS inputs to be the same size. Recycling is mainly useful for RHS inputs, where you might supply a size 1 input that will be recycled to the size of the LHS inputs.

NULL inputs are ignored.

.default	The value used when all of the LHS inputs return either FALSE or NA.
	.default must be size 1 or the same size as the common size computed from
	.default participates in the computation of the common type with the RHS inputs.
	NA values in the LHS conditions are treated like FALSE, meaning that the result at those locations will be assigned the .default value. To handle missing values in the conditions differently, you must explicitly catch them with another condition before they fall through to the .default. This typically involves some variation of is.na(x) ~ value tailored to your usage of case_when().
	If NULL, the default, a missing value will be used.
.ptype	An optional prototype declaring the desired output type. If supplied, this over- rides the common type of the RHS inputs.
.size	An optional size declaring the desired output size. If supplied, this overrides the common size computed from

## Value

A vector with the same size as the common size computed from the inputs in ... and the same type as the common type of the RHS inputs in ....

## See Also

case\_match()

```
x <- 1:70
case_when(
  x %% 35 == 0 ~ "fizz buzz",
 x %% 5 == 0 ~ "fizz",
  x %% 7 == 0 ~ "buzz",
  .default = as.character(x)
)
# Like an if statement, the arguments are evaluated in order, so you must
# proceed from the most specific to the most general. This won't work:
case_when(
 x %% 5 == 0 ~ "fizz",
  x %% 7 == 0 ~ "buzz",
 x %% 35 == 0 ~ "fizz buzz",
  .default = as.character(x)
)
# If none of the cases match and no `.default` is supplied, NA is used:
case_when(
 x %% 35 == 0 ~ "fizz buzz",
 x %% 5 == 0 ~ "fizz",
 x %% 7 == 0 ~ "buzz",
)
```

```
# Note that `NA` values on the LHS are treated like `FALSE` and will be
# assigned the `.default` value. You must handle them explicitly if you
# want to use a different value. The exact way to handle missing values is
# dependent on the set of LHS conditions you use.
x[2:4] <- NA_real_
case_when(
 x %% 35 == 0 ~ "fizz buzz",
 x %% 5 == 0 ~ "fizz",
 x %% 7 == 0 ~ "buzz",
 is.na(x) ~ "nope",
 .default = as.character(x)
)
# `case_when()` evaluates all RHS expressions, and then constructs its
# result by extracting the selected (via the LHS expressions) parts.
# In particular `NaN`s are produced in this case:
y \le seq(-2, 2, by = .5)
case_when(
 y \ge 0 \sim sqrt(y),
 .default = y
)
# `case_when()` is particularly useful inside `mutate()` when you want to
# create a new variable that relies on a complex combination of existing
# variables
starwars %>%
 select(name:mass, gender, species) %>%
 mutate(
   type = case_when(
     height > 200 | mass > 200 ~ "large",
     species == "Droid" ~ "robot",
      .default = "other"
   )
 )
# `case_when()` is not a tidy eval function. If you'd like to reuse
# the same patterns, extract the `case_when()` call in a normal
# function:
case_character_type <- function(height, mass, species) {</pre>
 case_when(
   height > 200 | mass > 200 ~ "large",
    species == "Droid" ~ "robot",
    .default = "other"
 )
}
case_character_type(150, 250, "Droid")
case_character_type(150, 150, "Droid")
# Such functions can be used inside `mutate()` as well:
starwars %>%
```

## coalesce

```
mutate(type = case_character_type(height, mass, species)) %>%
 pull(type)
# `case_when()` ignores `NULL` inputs. This is useful when you'd
# like to use a pattern only under certain conditions. Here we'll
# take advantage of the fact that `if` returns `NULL` when there is
# no `else` clause:
case_character_type <- function(height, mass, species, robots = TRUE) {</pre>
 case_when(
   height > 200 | mass > 200 ~ "large",
   if (robots) species == "Droid" ~ "robot",
    .default = "other"
 )
}
starwars %>%
 mutate(type = case_character_type(height, mass, species, robots = FALSE)) %>%
 pull(type)
```

```
coalesce
```

Find the first non-missing element

## Description

Given a set of vectors, coalesce() finds the first non-missing value at each position. It's inspired by the SQL COALESCE function which does the same thing for SQL NULLS.

## Usage

coalesce(..., .ptype = NULL, .size = NULL)

## Arguments

One or more vectors. These will be recycled against each other, and will be cast to their common type.
An optional prototype declaring the desired output type. If supplied, this over- rides the common type of the vectors in
An optional size declaring the desired output size. If supplied, this overrides the common size of the vectors in

## Value

A vector with the same type and size as the common type and common size of the vectors in ....

## See Also

na\_if() to replace specified values with an NA. tidyr::replace\_na() to replace NA with a value.

## Examples

```
# Use a single value to replace all missing values
x <- sample(c(1:5, NA, NA, NA))
coalesce(x, 0L)
# The equivalent to a missing value in a list is `NULL`
coalesce(list(1, 2, NULL), list(NA))
# Or generate a complete vector from partially missing pieces
y <- c(1, 2, NA, NA, 5)
z <- c(NA, NA, 3, 4, 5)
coalesce(y, z)
# Supply lists by splicing them into dots:
vecs <- list(
    c(1, 2, NA, NA, 5),
    c(NA, NA, 3, 4, 5)
)
coalesce(!!!vecs)
```

compute

Force computation of a database query

# Description

compute() stores results in a remote temporary table. collect() retrieves data into a local tibble. collapse() is slightly different: it doesn't force computation, but instead forces generation of the SQL query. This is sometimes needed to work around bugs in dplyr's SQL generation.

All functions preserve grouping and ordering.

## Usage

```
compute(x, ...)
collect(x, ...)
collapse(x, ...)
```

# Arguments

x	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	Arguments passed on to methods

#### consecutive\_id

#### Methods

These functions are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

- compute(): no methods found
- collect(): no methods found
- collapse(): no methods found

## See Also

copy\_to(), the opposite of collect(): it takes a local data frame and uploads it to the remote
source.

#### Examples

```
mtcars2 <- dbplyr::src_memdb() %>%
   copy_to(mtcars, name = "mtcars2-cc", overwrite = TRUE)
remote <- mtcars2 %>%
   filter(cyl == 8) %>%
   select(mpg:drat)
# Compute query and save in remote table
compute(remote)
# Compute query bring back to this session
collect(remote)
# Creates a fresh query based on the generated SQL
collapse(remote)
```

consecutive\_id *Generate a unique identifier for consecutive combinations* 

## Description

consecutive\_id() generates a unique identifier that increments every time a variable (or combination of variables) changes. Inspired by data.table::rleid().

#### Usage

consecutive\_id(...)

#### Arguments

• • •

Unnamed vectors. If multiple vectors are supplied, then they should have the same length.

## Value

A numeric vector the same length as the longest element of ....

#### Examples

```
consecutive_id(c(TRUE, TRUE, FALSE, FALSE, TRUE, FALSE, NA, NA))
consecutive_id(c(1, 1, 1, 2, 1, 1, 2, 2))
df <- data.frame(x = c(0, 0, 1, 0), y = c(2, 2, 2, 2))
df %>% group_by(x, y) %>% summarise(n = n())
df %>% group_by(id = consecutive_id(x, y), x, y) %>% summarise(n = n())
```

```
context
```

Information about the "current" group or variable

## Description

These functions return information about the "current" group or "current" variable, so only work inside specific contexts like summarise() and mutate().

- n() gives the current group size.
- cur\_group() gives the group keys, a tibble with one row and one column for each grouping variable.
- cur\_group\_id() gives a unique numeric identifier for the current group.
- cur\_group\_rows() gives the row indices for the current group.
- cur\_column() gives the name of the current column (in across() only).

See group\_data() for equivalent functions that return values for all groups.

See pick() for a way to select a subset of columns using tidyselect syntax while inside summarise() or mutate().

## Usage

n()

cur\_group()

cur\_group\_id()

cur\_group\_rows()

```
cur_column()
```

## copy\_to

# data.table

If you're familiar with data.table:

- cur\_group\_id() <-> .GRP
- cur\_group() <-> .BY
- cur\_group\_rows() <-> .I

See pick() for an equivalent to .SD.

# Examples

```
df <- tibble(
  g = sample(rep(letters[1:3], 1:3)),
  x = runif(6),
  y = runif(6)
)
gf <- df %>% group_by(g)
gf %>% summarise(n = n())
gf %>% mutate(id = cur_group_id())
gf %>% reframe(row = cur_group_rows())
gf %>% summarise(data = list(cur_group()))
gf %>% mutate(across(everything(), ~ paste(cur_column(), round(.x, 2))))
```

copy\_to

# Copy a local data frame to a remote src

# Description

This function uploads a local data frame into a remote data source, creating the table definition as needed. Wherever possible, the new object will be temporary, limited to the current connection to the source.

## Usage

```
copy_to(dest, df, name = deparse(substitute(df)), overwrite = FALSE, ...)
```

# Arguments

dest	remote data source
df	local data frame
name	name for new remote table.
overwrite	If TRUE, will overwrite an existing table with name name. If FALSE, will throw an error if name already exists.
	other parameters passed to methods.

#### Value

a tbl object in the remote source

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

## See Also

collect() for the opposite action; downloading remote data into a local dbl.

# Examples

```
## Not run:
iris2 <- dbplyr::src_memdb() %>% copy_to(iris, overwrite = TRUE)
iris2
```

## End(Not run)

count

Count the observations in each group

#### Description

count() lets you quickly count the unique values of one or more variables: df  $\gg$  count(a, b) is roughly equivalent to df  $\gg$  group\_by(a, b)  $\gg$  summarise(n = n()). count() is paired with tally(), a lower-level helper that is equivalent to df  $\gg$  summarise(n = n()). Supply wt to perform weighted counts, switching the summary from n = n() to n = sum(wt).

add\_count() and add\_tally() are equivalents to count() and tally() but use mutate() instead of summarise() so that they add a new column with group-wise counts.

# Usage

```
count(x, ..., wt = NULL, sort = FALSE, name = NULL)
## S3 method for class 'data.frame'
count(
    x,
    ...,
    wt = NULL,
    sort = FALSE,
    name = NULL,
    .drop = group_by_drop_default(x)
)
```

```
tally(x, wt = NULL, sort = FALSE, name = NULL)
add_count(x, ..., wt = NULL, sort = FALSE, name = NULL, .drop = deprecated())
```

```
add_tally(x, wt = NULL, sort = FALSE, name = NULL)
```

## Arguments

х	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr).
	<pre><data-masking> Variables to group by.</data-masking></pre>
wt	<pre><data-masking> Frequency weights. Can be NULL or a variable:</data-masking></pre>
	• If NULL (the default), counts the number of rows in each group.
	• If a variable, computes sum(wt) for each group.
sort	If TRUE, will show the largest groups at the top.
name	The name of the new column in the output.
	If omitted, it will default to n. If there's already a column called n, it will use nn. If there's a column called n and nn, it'll use nnn, and so on, adding ns until it gets a new name.
.drop	Handling of factor levels that don't appear in the data, passed on to group_by().
	For count(): if FALSE will include counts for empty groups (i.e. for levels of factors that don't exist in the data).
	[ <b>Deprecated</b> ] For add_count(): deprecated since it can't actually affect the output.

## Value

An object of the same type as .data. count() and add\_count() group transiently, so the output has the same groups as the input.

```
# count() is a convenient way to get a sense of the distribution of
# values in a dataset
starwars %>% count(species)
starwars %>% count(species, sort = TRUE)
starwars %>% count(sex, gender, sort = TRUE)
starwars %>% count(birth_decade = round(birth_year, -1))
# use the `wt` argument to perform a weighted count. This is useful
# when the data has already been aggregated once
df <- tribble(</pre>
  ~name,
            ~gender, ~runs,
  "Max",
            "male",
                        10,
  "Sandra", "female",
                         1,
  "Susan", "female",
                          4
)
```

```
# counts rows:
df %>% count(gender)
# counts runs:
df %>% count(gender, wt = runs)
# When factors are involved, `.drop = FALSE` can be used to retain factor
# levels that don't appear in the data
df2 <- tibble(
  id = 1:5,
  type = factor(c("a", "c", "a", NA, "a"), levels = c("a", "b", "c"))
)
df2 %>% count(type)
df2 %>% count(type, .drop = FALSE)
# Or, using `group_by()`:
df2 %>% group_by(type, .drop = FALSE) %>% count()
# tally() is a lower-level function that assumes you've done the grouping
starwars %>% tally()
starwars %>% group_by(species) %>% tally()
# both count() and tally() have add_ variants that work like
# mutate() instead of summarise
df %>% add_count(gender, wt = runs)
df %>% add_tally(wt = runs)
```

cross\_join Cross join

# Description

Cross joins match each row in x to every row in y, resulting in a data frame with nrow(x) \* nrow(y) rows.

Since cross joins result in all possible matches between x and y, they technically serve as the basis for all mutating joins, which can generally be thought of as cross joins followed by a filter. In practice, a more specialized procedure is used for better performance.

## Usage

```
cross_join(x, y, ..., copy = FALSE, suffix = c(".x", ".y"))
```

## Arguments

х, у	A pair of data frames, data frame extensions (e.g. a tibble), or lazy data frames (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	Other parameters passed onto methods.
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.

#### cumall

suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added
	to the output to disambiguate them. Should be a character vector of length 2.

## Value

An object of the same type as x (including the same groups). The output has the following properties:

- There are nrow(x) \* nrow(y) rows returned.
- Output columns include all columns from both x and y. Column name collisions are resolved using suffix.
- The order of the rows and columns of x is preserved as much as possible.

# Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

## See Also

Other joins: filter-joins, mutate-joins, nest\_join()

#### Examples

```
# Cross joins match each row in `x` to every row in `y`.
# Data within the columns is not used in the matching process.
cross_join(band_instruments, band_members)
# Control the suffix added to variables duplicated in
```

```
# `x` and `y` with `suffix`.
cross_join(band_instruments, band_members, suffix = c("", "_y"))
```

cumall

Cumulativate versions of any, all, and mean

## Description

dplyr provides cumal1(), cumany(), and cummean() to complete R's set of cumulative functions.

## Usage

cumany(x)

cumall(x)

cummean(x)

## Arguments

Х

For cumall() and cumany(), a logical vector; for cummean() an integer or numeric vector.

# Value

A vector the same length as x.

# **Cumulative logical functions**

These are particularly useful in conjunction with filter():

- cumall(x): all cases until the first FALSE.
- cumall(!x): all cases until the first TRUE.
- cumany(x): all cases after the first TRUE.
- cumany(!x): all cases after the first FALSE.

```
# `cummean()` returns a numeric/integer vector of the same length
# as the input vector.
x <- c(1, 3, 5, 2, 2)
cummean(x)
cumsum(x) / seq_along(x)
# `cumall()` and `cumany()` return logicals
cumall(x < 5)
cumany(x == 3)
# cumall() vs. cumany()
df <- data.frame(</pre>
  date = as.Date("2020-01-01") + 0:6,
  balance = c(100, 50, 25, -25, -50, 30, 120)
)
# all rows after first overdraft
df %>% filter(cumany(balance < 0))</pre>
# all rows until first overdraft
df %>% filter(cumall(!(balance < 0)))</pre>
```

desc

## Description

c\_across() is designed to work with rowwise() to make it easy to perform row-wise aggregations. It has two differences from c():

- It uses tidy select semantics so you can easily select multiple variables. See vignette("rowwise") for more details.
- It uses vctrs::vec\_c() in order to give safer outputs.

## Usage

```
c_across(cols)
```

# Arguments

cols

<tidy-select> Columns to transform. You can't select grouping columns because they are already automatically handled by the verb (i.e. summarise() or mutate()).

# See Also

across() for a function that returns a tibble.

## Examples

```
df <- tibble(id = 1:4, w = runif(4), x = runif(4), y = runif(4), z = runif(4))
df %>%
rowwise() %>%
mutate(
    sum = sum(c_across(w:z)),
    sd = sd(c_across(w:z))
)
```

```
desc
```

Descending order

## Description

Transform a vector into a format that will be sorted in descending order. This is useful within arrange().

## Usage

desc(x)

## Arguments

x vector to transform

## distinct

## Examples

```
desc(1:10)
desc(factor(letters))
first_day <- seq(as.Date("1910/1/1"), as.Date("1920/1/1"), "years")
desc(first_day)
starwars %>% arrange(desc(mass))
```

```
distinct
```

Keep distinct/unique rows

# Description

Keep only unique/distinct rows from a data frame. This is similar to unique.data.frame() but considerably faster.

## Usage

distinct(.data, ..., .keep\_all = FALSE)

## Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g.
	from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<pre><data-masking> Optional variables to use when determining uniqueness. If</data-masking></pre>
	there are multiple rows for a given combination of inputs, only the first row will
	be preserved. If omitted, will use all variables in the data frame.
.keep_all	If TRUE, keep all variables in .data. If a combination of is not distinct, this
	keeps the first row of values.

## Value

An object of the same type as . data. The output has the following properties:

- Rows are a subset of the input but appear in the same order.
- Columns are not modified if ... is empty or .keep\_all is TRUE. Otherwise, distinct() first calls mutate() to create new columns.
- Groups are not modified.
- Data frame attributes are preserved.

# Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

## dplyr\_by

## Examples

```
df <- tibble(</pre>
 x = sample(10, 100, rep = TRUE),
 y = sample(10, 100, rep = TRUE)
)
nrow(df)
nrow(distinct(df))
nrow(distinct(df, x, y))
distinct(df, x)
distinct(df, y)
# You can choose to keep all other variables as well
distinct(df, x, .keep_all = TRUE)
distinct(df, y, .keep_all = TRUE)
# You can also use distinct on computed variables
distinct(df, diff = abs(x - y))
# Use `pick()` to select columns with tidy-select
distinct(starwars, pick(contains("color")))
# Grouping -----
df <- tibble(</pre>
 g = c(1, 1, 2, 2, 2),
 x = c(1, 1, 2, 1, 2),
 y = c(3, 2, 1, 3, 1)
)
df <- df %>% group_by(g)
# With grouped data frames, distinctness is computed within each group
df %>% distinct(x)
# When `...` are omitted, `distinct()` still computes distinctness using
# all variables in the data frame
df %>% distinct()
```

dplyr\_by

# Per-operation grouping with .by/by

# Description

There are two ways to group in dplyr:

- Persistent grouping with group\_by()
- Per-operation grouping with .by/by

This help page is dedicated to explaining where and why you might want to use the latter.

Depending on the dplyr verb, the per-operation grouping argument may be named .by or by. The *Supported verbs* section below outlines this on a case-by-case basis. The remainder of this page will refer to .by for simplicity.

Grouping radically affects the computation of the dplyr verb you use it with, and one of the goals of .by is to allow you to place that grouping specification alongside the code that actually uses it. As an added benefit, with .by you no longer need to remember to ungroup() after summarise(), and summarise() won't ever message you about how it's handling the groups!

This idea comes from data.table, which allows you to specify by alongside modifications in j, like: dt[, .(x = mean(x)), by = g].

#### Supported verbs:

- mutate(.by = )
- summarise(.by = )
- reframe(.by = )
- filter(.by = )
- slice(.by = )
- slice\_head(by = ) and slice\_tail(by = )
- slice\_min(by = ) and slice\_max(by = )
- slice\_sample(by = )

Note that some dplyr verbs use by while others use .by. This is a purely technical difference.

Differences between .by and group\_by():

.by	group_by()
Grouping only affects a single verb	Grouping is persistent across multiple verbs
Selects variables with tidy-select	Computes expressions with data-masking
Summaries use existing order of group keys	Summaries sort group keys in ascending order

## Using .by:

Let's take a look at the two grouping approaches using this expenses data set, which tracks costs accumulated across various ids and regions:

```
expenses <- tibble(</pre>
  id = c(1, 2, 1, 3, 1, 2, 3),
  region = c("A", "A", "A", "B", "B", "A", "A"),
  cost = c(25, 20, 19, 12, 9, 6, 6)
)
expenses
#> # A tibble: 7 x 3
#>
        id region cost
#>
     <dbl> <chr> <dbl>
#> 1
        1 A
                     25
#> 2
         2 A
                     20
#> 3
         1 A
                     19
#> 4
         3 B
                     12
```

dplyr\_by

#> 5	5 1	B 9	)
#> 6	5 2	Α 6	5
#> 7	3	A 6	5

Imagine that you wanted to compute the average cost per region. You'd probably write something like this:

```
expenses %>%
group_by(region) %>%
summarise(cost = mean(cost))
#> # A tibble: 2 x 2
#> region cost
#> <chr> <dbl>
#> 1 A 15.2
#> 2 B 10.5
```

Instead, you can now specify the grouping *inline* within the verb:

```
expenses %>%
  summarise(cost = mean(cost), .by = region)
#> # A tibble: 2 x 2
#> region cost
#> <chr> <dbl>
#> 1 A 15.2
#> 2 B 10.5
```

. by applies to a single operation, meaning that since expenses was an ungrouped data frame, the result after applying .by will also always be an ungrouped data frame, regardless of the number of grouping columns.

```
expenses %>%
```

```
summarise(cost = mean(cost), .by = c(id, region))
#> # A tibble: 5 x 3
#>
        id region cost
#>
    <dbl> <chr> <dbl>
#> 1
        1 A
                     22
#> 2
                     13
         2 A
#> 3
         3 B
                     12
#> 4
         1 B
                      9
#> 5
         3 A
                      6
```

Compare that with group\_by() %>% summarise(), where summarise() generally peels off 1 layer of grouping by default, typically with a message that it is doing so:

```
expenses %>%
 group_by(id, region) %>%
 summarise(cost = mean(cost))
#> `summarise()` has grouped output by 'id'. You can override using the `.groups`
#> argument.
#> # A tibble: 5 x 3
#> # Groups: id [3]
#> id region cost
```

#>		<dbl></dbl>	<chr></chr>	<dbl></dbl>
#>	1	1	А	22
#>	2	1	В	9
#>	3	2	А	13
#>	4	3	А	6
#>	5	3	В	12

Because .by grouping applies to a single operation, you don't need to worry about ungrouping, and it never needs to emit a message to remind you what it is doing with the groups.

Note that with .by we specified multiple columns to group by using the tidy-select syntax c(id, region). If you have a character vector of column names you'd like to group by, you can do so with .by = all\_of(my\_cols). It will group by the columns in the order they were provided.

To prevent surprising results, you can't use . by on an existing grouped data frame:

```
expenses %>%
  group_by(id) %>%
  summarise(cost = mean(cost), .by = c(id, region))
#> Error in `summarise()`:
#> ! Can't supply `.by` when `.data` is a grouped data frame.
```

So far we've focused on the usage of .by with summarise(), but .by works with a number of other dplyr verbs. For example, you could append the mean cost per region onto the original data frame as a new column rather than computing a summary:

```
expenses %>%
  mutate(cost_by_region = mean(cost), .by = region)
#> # A tibble: 7 x 4
#>
        id region cost cost_by_region
#>
     <dbl> <chr> <dbl>
                                  <dbl>
#> 1
         1 A
                     25
                                   15.2
#> 2
         2 A
                      20
                                   15.2
#> 3
         1 A
                     19
                                   15.2
#> 4
         3 B
                     12
                                   10.5
#> 5
         1 B
                      9
                                   10.5
                      6
#> 6
         2 A
                                   15.2
#> 7
         3 A
                      6
                                   15.2
```

Or you could slice out the maximum cost per combination of id and region:

```
# Note that the argument is named `by` in `slice_max()`
expenses %>%
  slice_max(cost, n = 1, by = c(id, region))
#> # A tibble: 5 x 3
#>
        id region cost
#>
     <dbl> <chr> <dbl>
#> 1
         1 A
                     25
#> 2
         2 A
                     20
#> 3
         3 B
                     12
#> 4
         1 B
                      9
#> 5
                      6
         3 A
```

## dplyr\_by

## **Result ordering:**

When used with .by, summarise(), reframe(), and slice() all maintain the ordering of the existing data. This is different from group\_by(), which has always sorted the group keys in ascending order.

```
df <- tibble(</pre>
  month = c("jan", "jan", "feb", "feb", "mar"),
  temp = c(20, 25, 18, 20, 40)
)
# Uses ordering by "first appearance" in the original data
df %>%
  summarise(average_temp = mean(temp), .by = month)
#> # A tibble: 3 x 2
     month average_temp
#>
#>
     <chr>
                 <dbl>
#> 1 jan
                   22.5
#> 2 feb
                   19
#> 3 mar
                   40
# Sorts in ascending order
df %>%
  group_by(month) %>%
  summarise(average_temp = mean(temp))
#> # A tibble: 3 x 2
#>
     month average_temp
#>
     <chr>
                  <dbl>
#> 1 feb
                   19
#> 2 jan
                   22.5
#> 3 mar
                   40
```

If you need sorted group keys, we recommend that you explicitly use arrange() either before or after the call to summarise(), reframe(), or slice(). This also gives you full access to all of arrange()'s features, such as desc() and the .locale argument.

#### Verbs without .by support:

If a dplyr verb doesn't support .by, then that typically means that the verb isn't inherently affected by grouping. For example, pull() and rename() don't support .by, because specifying columns to group by would not affect their implementations.

That said, there are a few exceptions to this where sometimes a dplyr verb doesn't support .by, but *does* have special support for grouped data frames created by group\_by(). This is typically because the verbs are required to retain the grouping columns, for example:

- select() always retains grouping columns, with a message if any aren't specified in the select() call.
- distinct() and count() place unspecified grouping columns at the front of the data frame before computing their results.

• arrange() has a .by\_group argument to optionally order by grouping columns first.

If group\_by() didn't exist, then these verbs would not have special support for grouped data frames.

explain

#### Description

This is a generic function which gives more details about an object than print(), and is more focused on human readable output than str().

## Usage

explain(x, ...)

show\_query(x, ...)

## Arguments

х	An object to explain
	Other parameters possibly used by generic

## Value

The first argument, invisibly.

## Databases

Explaining a tbl\_sql will run the SQL EXPLAIN command which will describe the query plan. This requires a little bit of knowledge about how EXPLAIN works for your database, but is very useful for diagnosing performance problems.

```
lahman_s <- dbplyr::lahman_sqlite()
batting <- tbl(lahman_s, "Batting")
batting %>% show_query()
batting %>% explain()
# The batting database has indices on all ID variables:
# SQLite automatically picks the most restrictive index
batting %>% filter(lgID == "NL" & yearID == 2000L) %>% explain()
# OR's will use multiple indexes
batting %>% filter(lgID == "NL" | yearID == 2000) %>% explain()
# Joins will use indexes in both tables
teams <- tbl(lahman_s, "Teams")
batting %>% left_join(teams, c("yearID", "teamID")) %>% explain()
```
filter

### Description

The filter() function is used to subset a data frame, retaining all rows that satisfy your conditions. To be retained, the row must produce a value of TRUE for all conditions. Note that when a condition evaluates to NA the row will be dropped, unlike base subsetting with [.

#### Usage

filter(.data, ..., .by = NULL, .preserve = FALSE)

### Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<data-masking> Expressions that return a logical value, and are defined in terms of the variables in .data. If multiple expressions are included, they are combined with the &amp; operator. Only rows for which all conditions evaluate to TRUE are kept.</data-masking>
.by	[Experimental]
	<tidy-select> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to group_by(). For details and examples, see ?dplyr_by.</tidy-select>
.preserve	Relevant when the .data input is grouped. If .preserve = FALSE (the default), the grouping structure is recalculated based on the resulting data, otherwise the grouping is kept as is.

# Details

The filter() function is used to subset the rows of .data, applying the expressions in ... to the column values to determine which rows should be retained. It can be applied to both grouped and ungrouped data (see group\_by() and ungroup()). However, dplyr is not yet smart enough to optimise the filtering operation on grouped datasets that do not need grouped calculations. For this reason, filtering is often considerably faster on ungrouped data.

### Value

An object of the same type as . data. The output has the following properties:

- Rows are a subset of the input, but appear in the same order.
- Columns are not modified.
- The number of groups may be reduced (if .preserve is not TRUE).
- Data frame attributes are preserved.

#### **Useful filter functions**

There are many functions and operators that are useful when constructing the expressions used to filter the data:

- ==, >, >= etc
- &, |, !, xor()
- is.na()
- between(), near()

#### Grouped tibbles

Because filtering expressions are computed within groups, they may yield different results on grouped tibbles. This will be the case as soon as an aggregating, lagging, or ranking function is involved. Compare this ungrouped filtering:

```
starwars %>% filter(mass > mean(mass, na.rm = TRUE))
```

With the grouped equivalent:

```
starwars %>% group_by(gender) %>% filter(mass > mean(mass, na.rm = TRUE))
```

In the ungrouped version, filter() compares the value of mass in each row to the global average (taken over the whole data set), keeping only the rows with mass greater than this global average. In contrast, the grouped version calculates the average mass separately for each gender group, and keeps rows with mass greater than the relevant within-gender average.

#### Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

#### See Also

Other single table verbs: arrange(), mutate(), reframe(), rename(), select(), slice(), summarise()

#### Examples

```
# Filtering by one criterion
filter(starwars, species == "Human")
filter(starwars, mass > 1000)
# Filtering by multiple criteria within a single logical expression
filter(starwars, hair_color == "none" & eye_color == "black")
filter(starwars, hair_color == "none" | eye_color == "black")
```

# When multiple expressions are used, they are combined using &

#### filter-joins

```
filter(starwars, hair_color == "none", eye_color == "black")
# The filtering operation may yield different results on grouped
# tibbles because the expressions are computed within groups.
#
# The following filters rows where `mass` is greater than the
# global average:
starwars %>% filter(mass > mean(mass, na.rm = TRUE))
# Whereas this keeps rows with `mass` greater than the gender
# average:
starwars %>% group_by(gender) %>% filter(mass > mean(mass, na.rm = TRUE))
# To refer to column names that are stored as strings, use the `.data` pronoun:
vars <- c("mass", "height")</pre>
cond <- c(80, 150)
starwars %>%
 filter(
    .data[[vars[[1]]]] > cond[[1]],
    .data[[vars[[2]]] > cond[[2]]
 )
# Learn more in ?rlang::args_data_masking
```

filter-joins Filtering joins

### Description

Filtering joins filter rows from x based on the presence or absence of matches in y:

- semi\_join() return all rows from x with a match in y.
- anti\_join() return all rows from x with**out** a match in y.

### Usage

```
semi_join(x, y, by = NULL, copy = FALSE, ...)
## S3 method for class 'data.frame'
semi_join(x, y, by = NULL, copy = FALSE, ..., na_matches = c("na", "never"))
anti_join(x, y, by = NULL, copy = FALSE, ...)
## S3 method for class 'data.frame'
anti_join(x, y, by = NULL, copy = FALSE, ..., na_matches = c("na", "never"))
```

## Arguments

х, у	A pair of data frames, data frame extensions (e.g. a tibble), or lazy data frames (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
by	A join specification created with join_by(), or a character vector of variables to join by.
	If NULL, the default, *_join() will perform a natural join, using all variables in common across x and y. A message lists the variables so that you can check they're correct; suppress the message by supplying by explicitly.
	To join on different variables between x and y, use a join_by() specification. For example, join_by(a == b) will match x\$a to y\$b.
	To join by multiple variables, use a join_by() specification with multiple expressions. For example, join_by(a == b, c == d) will match x\$a to y\$b and x\$c to y\$d. If the column names are the same between x and y, you can shorten this by listing only the variable names, like join_by(a, c).
	join_by() can also be used to perform inequality, rolling, and overlap joins. See the documentation at ?join_by for details on these types of joins.
	For simple equality joins, you can alternatively specify a character vector of variable names to join by. For example, by = $c("a", "b")$ joins x\$a to y\$a and x\$b to y\$b. If variable names differ between x and y, use a named character vector like by = $c("x_a" = "y_a", "x_b" = "y_b")$ .
	To perform a cross-join, generating all combinations of x and y, see cross_join().
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.
	Other parameters passed onto methods.
na_matches	Should two NA or two NaN values match?
	<ul> <li>"na", the default, treats two NA or two NaN values as equal, like %in%, match(), and merge().</li> </ul>
	<ul> <li>"never" treats two NA or two NaN values as different, and will never match them together or to any other values. This is similar to joins for database sources and to base::merge(incomparables = NA).</li> </ul>

## Value

An object of the same type as x. The output has the following properties:

- Rows are a subset of the input, but appear in the same order.
- Columns are not modified.
- Data frame attributes are preserved.
- Groups are taken from x. The number of groups may be reduced.

## Methods

These function are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

#### glimpse

Methods available in currently loaded packages:

- semi\_join(): no methods found.
- anti\_join(): no methods found.

### See Also

Other joins: cross\_join(), mutate-joins, nest\_join()

#### Examples

```
# "Filtering" joins keep cases from the LHS
band_members %>% semi_join(band_instruments)
band_members %>% anti_join(band_instruments)
```

```
# To suppress the message about joining variables, supply `by`
band_members %>% semi_join(band_instruments, by = join_by(name))
# This is good practice in production code
```

glimpse

Get a glimpse of your data

#### Description

glimpse() is like a transposed version of print(): columns run down the page, and data runs across. This makes it possible to see every column in a data frame. It's a little like str() applied to a data frame but it tries to show you as much data as possible. (And it always shows the underlying data, even when applied to a remote data source.)

glimpse() is provided by the pillar package, and re-exported by dplyr. See pillar::glimpse() for more details.

## Value

x original x is (invisibly) returned, allowing glimpse() to be used within a data pipeline.

```
glimpse(mtcars)
```

```
# Note that original x is (invisibly) returned, allowing `glimpse()` to be
# used within a pipeline.
mtcars %>%
  glimpse() %>%
  select(1:3)
glimpse(starwars)
```

group\_by

### Description

Most data operations are done on groups defined by variables. group\_by() takes an existing tbl and converts it into a grouped tbl where operations are performed "by group". ungroup() removes grouping.

### Usage

```
group_by(.data, ..., .add = FALSE, .drop = group_by_drop_default(.data))
```

ungroup(x, ...)

### Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	In group_by(), variables or computations to group by. Computations are always done on the ungrouped data frame. To perform computations on the grouped data, you need to use a separate mutate() step before the group_by(). Computations are not allowed in nest_by(). In ungroup(), variables to remove from the grouping.
.add	When FALSE, the default, group_by() will override existing groups. To add to the existing groups, use .add = TRUE.
	This argument was previously called add, but that prevented creating a new grouping variable called add, and conflicts with our naming conventions.
.drop	Drop groups formed by factor levels that don't appear in the data? The default is TRUE except when .data has been previously grouped with .drop = FALSE. See group_by_drop_default() for details.
x	A tbl()

### Value

A grouped data frame with class grouped\_df, unless the combination of . . . and add yields a empty set of grouping columns, in which case a tibble will be returned.

#### Methods

These function are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

- group\_by(): no methods found.
- ungroup(): no methods found.

group\_by

### Ordering

Currently, group\_by() internally orders the groups in ascending order. This results in ordered output from functions that aggregate groups, such as summarise().

When used as grouping columns, character vectors are ordered in the C locale for performance and reproducibility across R sessions. If the resulting ordering of your grouped operation matters and is dependent on the locale, you should follow up the grouped operation with an explicit call to arrange() and set the .locale argument. For example:

data %>%

group\_by(chr) %>%
summarise(avg = mean(x)) %>%
arrange(chr, .locale = "en")

This is often useful as a preliminary step before generating content intended for humans, such as an HTML table.

### Legacy behavior:

Prior to dplyr 1.1.0, character vector grouping columns were ordered in the system locale. If you need to temporarily revert to this behavior, you can set the global option dplyr.legacy\_locale to TRUE, but this should be used sparingly and you should expect this option to be removed in a future version of dplyr. It is better to update existing code to explicitly call arrange(.locale = ) instead. Note that setting dplyr.legacy\_locale will also force calls to arrange() to use the system locale.

### See Also

Other grouping functions: group\_map(), group\_nest(), group\_split(), group\_trim()

### Examples

```
by_cyl <- mtcars %>% group_by(cyl)
```

```
# grouping doesn't change how the data looks (apart from listing
# how it's grouped):
by_cyl
# It changes how it acts with the other dplyr verbs:
by_cyl %>% summarise(
  disp = mean(disp),
  hp = mean(hp)
)
by_cyl %>% filter(disp == max(disp))
# Each call to summarise() removes a layer of grouping
by_vs_am <- mtcars %>% group_by(vs, am)
by_vs <- by_vs_am %>% summarise(n = n())
by_vs %>% summarise(n = sum(n))
```

# To removing grouping, use ungroup

```
by_vs %>%
 ungroup() %>%
  summarise(n = sum(n))
# By default, group_by() overrides existing grouping
by_cyl %>%
  group_by(vs, am) %>%
  group_vars()
# Use add = TRUE to instead append
by_cyl %>%
  group_by(vs, am, .add = TRUE) %>%
  group_vars()
# You can group by expressions: this is a short-hand
# for a mutate() followed by a group_by()
mtcars %>%
  group_by(vsam = vs + am)
# The implicit mutate() step is always performed on the
# ungrouped data. Here we get 3 groups:
mtcars %>%
  group_by(vs) %>%
  group_by(hp_cut = cut(hp, 3))
# If you want it to be performed by groups,
# you have to use an explicit mutate() call.
# Here we get 3 groups per value of vs
mtcars %>%
  group_by(vs) %>%
  mutate(hp_cut = cut(hp, 3)) %>%
  group_by(hp_cut)
# when factors are involved and .drop = FALSE, groups can be empty
tbl <- tibble(</pre>
  x = 1:10,
  y = factor(rep(c("a", "c"), each = 5), levels = c("a", "b", "c"))
)
tbl %>%
  group_by(y, .drop = FALSE) %>%
  group_rows()
```

group\_cols

Select grouping variables

#### Description

This selection helpers matches grouping variables. It can be used in select() or vars() selections.

44

#### group\_map

### Usage

group\_cols(vars = NULL, data = NULL)

### Arguments

vars	Deprecated; please use data instead.
data	For advanced use only. The default NULL automatically finds the "current" data
	frames.

## See Also

groups() and group\_vars() for retrieving the grouping variables outside selection contexts.

## Examples

```
gdf <- iris %>% group_by(Species)
gdf %>% select(group_cols())
# Remove the grouping variables from mutate selections:
gdf %>% mutate_at(vars(-group_cols()), `/`, 100)
# -> No longer necessary with across()
gdf %>% mutate(across(everything(), ~ . / 100))
```

group\_map

Apply a function to each group

### Description

#### [Experimental]

group\_map(), group\_modify() and group\_walk() are purrr-style functions that can be used to iterate on grouped tibbles.

### Usage

```
group_map(.data, .f, ..., .keep = FALSE)
group_modify(.data, .f, ..., .keep = FALSE)
group_walk(.data, .f, ..., .keep = FALSE)
```

#### Arguments

A grouped tibble
A function or formula to apply to each group.
If a <b>function</b> , it is used as is. It should have at least 2 formal arguments.
If a <b>formula</b> , e.g. ~ head(.x), it is converted to a function.
In the formula, you can use

	• . or .x to refer to the subset of rows of .tbl for the given group
	• .y to refer to the key, a one row tibble with one column per grouping vari- able that identifies the group
	Additional arguments passed on to .f
.keep	are the grouping variables kept in .x

## Details

Use group\_modify() when summarize() is too limited, in terms of what you need to do and return for each group\_modify() is good for "data frame in, data frame out". If that is too limited, you need to use a nested or split workflow. group\_modify() is an evolution of do(), if you have used that before.

Each conceptual group of the data frame is exposed to the function . f with two pieces of information:

- The subset of the data for the group, exposed as .x.
- The key, a tibble with exactly one row and columns for each grouping variable, exposed as .y.

For completeness, group\_modify(), group\_map and group\_walk() also work on ungrouped data frames, in that case the function is applied to the entire data frame (exposed as .x), and .y is a one row tibble with no column, consistently with group\_keys().

### Value

- group\_modify() returns a grouped tibble. In that case . f must return a data frame.
- group\_map() returns a list of results from calling . f on each group.
- group\_walk() calls . f for side effects and returns the input . tbl, invisibly.

#### See Also

Other grouping functions: group\_by(), group\_nest(), group\_split(), group\_trim()

```
# return a list
mtcars %>%
group_by(cyl) %>%
group_map(~ head(.x, 2L))
# return a tibble grouped by `cyl` with 2 rows per group
# the grouping data is recalculated
mtcars %>%
group_by(cyl) %>%
group_by(cyl) %>%
group_modify(~ head(.x, 2L))
# a list of tibbles
iris %>%
group_by(Species) %>%
```

```
group_map(~ broom::tidy(lm(Petal.Length ~ Sepal.Length, data = .x)))
# a restructured grouped tibble
iris %>%
 group_by(Species) %>%
 group_modify(~ broom::tidy(lm(Petal.Length ~ Sepal.Length, data = .x)))
# a list of vectors
iris %>%
 group_by(Species) %>%
 group_map(~ quantile(.x$Petal.Length, probs = c(0.25, 0.5, 0.75)))
# to use group_modify() the lambda must return a data frame
iris %>%
 group_by(Species) %>%
 group_modify(~ {
     quantile(.x$Petal.Length, probs = c(0.25, 0.5, 0.75)) %>%
     tibble::enframe(name = "prob", value = "quantile")
 })
iris %>%
 group_by(Species) %>%
 group_modify(~ {
    .x %>%
     purrr::map_dfc(fivenum) %>%
     mutate(nms = c("min", "Q1", "median", "Q3", "max"))
 })
# group_walk() is for side effects
dir.create(temp <- tempfile())</pre>
iris %>%
 group_by(Species) %>%
 group_walk(~ write.csv(.x, file = file.path(temp, paste0(.y$Species, ".csv"))))
list.files(temp, pattern = "csv$")
unlink(temp, recursive = TRUE)
# group_modify() and ungrouped data frames
mtcars %>%
 group_modify(~ head(.x, 2L))
```

group\_trim

```
Trim grouping structure
```

#### Description

**[Experimental]** Drop unused levels of all factors that are used as grouping variables, then recalculates the grouping structure.

group\_trim() is particularly useful after a filter() that is intended to select a subset of groups.

### Usage

group\_trim(.tbl, .drop = group\_by\_drop\_default(.tbl))

## Arguments

.tbl	A grouped data frame
.drop	See group_by()

### Value

A grouped data frame

# See Also

Other grouping functions: group\_by(), group\_map(), group\_nest(), group\_split()

#### Examples

```
iris %>%
group_by(Species) %>%
filter(Species == "setosa", .preserve = TRUE) %>%
group_trim()
```

ident

Flag a character vector as SQL identifiers

### Description

ident() takes unquoted strings and flags them as identifiers. ident\_q() assumes its input has already been quoted, and ensures it does not get quoted again. This is currently used only for schema.table.

### Usage

ident(...)

#### Arguments

... A character vector, or name-value pairs

## Examples

```
# Identifiers are escaped with "
```

ident("x")

if\_else

### Description

if\_else() is a vectorized if-else. Compared to the base R equivalent, ifelse(), this function allows you to handle missing values in the condition with missing and always takes true, false, and missing into account when determining what the output type should be.

#### Usage

```
if_else(condition, true, false, missing = NULL, ..., ptype = NULL, size = NULL)
```

## Arguments

condition	A logical vector
true, false	Vectors to use for TRUE and FALSE values of condition.
	Both true and false will be recycled to the size of condition.
	true, false, and missing (if used) will be cast to their common type.
missing	If not NULL, will be used as the value for NA values of condition. Follows the same size and type rules as true and false.
	These dots are for future extensions and must be empty.
ptype	An optional prototype declaring the desired output type. If supplied, this over- rides the common type of true, false, and missing.
size	An optional size declaring the desired output size. If supplied, this overrides the size of condition.

## Value

A vector with the same size as condition and the same type as the common type of true, false, and missing.

Where condition is TRUE, the matching values from true, where it is FALSE, the matching values from false, and where it is NA, the matching values from missing, if provided, otherwise a missing value will be used.

```
x <- c(-5:5, NA)
if_else(x < 0, NA, x)
# Explicitly handle `NA` values in the `condition` with `missing`
if_else(x < 0, "negative", "positive", missing = "missing")
# Unlike `ifelse()`, `if_else()` preserves types
x <- factor(sample(letters[1:5], 10, replace = TRUE))
ifelse(x %in% c("a", "b", "c"), x, NA)</pre>
```

join\_by

```
if_else(x %in% c("a", "b", "c"), x, NA)
# `if_else()` is often useful for creating new columns inside of `mutate()`
starwars %>%
    mutate(category = if_else(height < 100, "short", "tall"), .keep = "used")</pre>
```

join\_by

Join specifications

#### Description

join\_by() constructs a specification that describes how to join two tables using a small domain specific language. The result can be supplied as the by argument to any of the join functions (such as left\_join()).

### Usage

join\_by(...)

#### Arguments

• • •

Expressions specifying the join.

Each expression should consist of one of the following:

- Equality condition: ==
- Inequality conditions: >=, >, <=, or <
- Rolling helper: closest()
- Overlap helpers: between(), within(), or overlaps()

Other expressions are not supported. If you need to perform a join on a computed variable, e.g. join\_by(sales\_date - 40 >= promo\_date), you'll need to precompute and store it in a separate column.

Column names should be specified as quoted or unquoted names. By default, the name on the left-hand side of a join condition refers to the left-hand table, unless overridden by explicitly prefixing the column name with either x or y. If a single column name is provided without any join conditions, it is interpreted as if that column name was duplicated on each side of ==, i.e. x is interpreted as x == x.

## Join types

The following types of joins are supported by dplyr:

- · Equality joins
- · Inequality joins
- · Rolling joins
- Overlap joins

50

join\_by

• Cross joins

Equality, inequality, rolling, and overlap joins are discussed in more detail below. Cross joins are implemented through cross\_join().

## **Equality joins:**

Equality joins require keys to be equal between one or more pairs of columns, and are the most common type of join. To construct an equality join using join\_by(), supply two column names to join with separated by ==. Alternatively, supplying a single name will be interpreted as an equality join between two columns of the same name. For example, join\_by(x) is equivalent to join\_by(x == x).

## **Inequality joins:**

Inequality joins match on an inequality, such as >, >=, <, or <=, and are common in time series analysis and genomics. To construct an inequality join using join\_by(), supply two column names separated by one of the above mentioned inequalities.

Note that inequality joins will match a single row in x to a potentially large number of rows in y. Be extra careful when constructing inequality join specifications!

### **Rolling joins:**

Rolling joins are a variant of inequality joins that limit the results returned from an inequality join condition. They are useful for "rolling" the closest match forward/backwards when there isn't an exact match. To construct a rolling join, wrap an inequality with closest().

- closest(expr)
  - expr must be an inequality involving one of: >, >=, <, or <=.

For example,  $closest(x \ge y)$  is interpreted as: For each value in x, find the closest value in y that is less than or equal to that x value.

closest() will always use the left-hand table (x) as the primary table, and the right-hand table (y) as the one to find the closest match in, regardless of how the inequality is specified. For example,  $closest(ya \ge xb)$  will always be interpreted as  $closest(xb \le ya)$ .

### **Overlap joins:**

Overlap joins are a special case of inequality joins involving one or two columns from the lefthand table *overlapping* a range defined by two columns from the right-hand table. There are three helpers that join\_by() recognizes to assist with constructing overlap joins, all of which can be constructed from simpler inequalities.

between(x, y\_lower, y\_upper, ..., bounds = "[]")
 For each value in x, this finds everywhere that value falls between [y\_lower, y\_upper].
 Equivalent to x >= y\_lower, x <= y\_upper by default.</li>

bounds can be one of "[]", "[)", "(]", or "()" to alter the inclusiveness of the lower and upper bounds. This changes whether  $\geq$  or  $\geq$  and  $\leq$  or  $\leq$  are used to build the inequalities shown above.

Dots are for future extensions and must be empty.

within(x\_lower, x\_upper, y\_lower, y\_upper)

For each range in [x\_lower, x\_upper], this finds everywhere that range falls completely within [y\_lower, y\_upper]. Equivalent to x\_lower >= y\_lower, x\_upper <= y\_upper. The inequalities used to build within() are the same regardless of the inclusiveness of the supplied ranges.

overlaps(x\_lower, x\_upper, y\_lower, y\_upper, ..., bounds = "[]")
For each range in [x\_lower, x\_upper], this finds everywhere that range overlaps [y\_lower, y\_upper]
in any capacity. Equivalent to x\_lower <= y\_upper, x\_upper >= y\_lower by default.
bounds can be one of "[]", "[)", "(]", or "()" to alter the inclusiveness of the lower and
upper bounds. "[]" uses <= and >=, but the 3 other options use < and > and generate the exact
same inequalities.

Dots are for future extensions and must be empty.

These conditions assume that the ranges are well-formed and non-empty, i.e. x\_lower <= x\_upper when bounds are treated as "[]", and x\_lower < x\_upper otherwise.

#### **Column referencing**

When specifying join conditions,  $join_by()$  assumes that column names on the left-hand side of the condition refer to the left-hand table (x), and names on the right-hand side of the condition refer to the right-hand table (y). Occasionally, it is clearer to be able to specify a right-hand table name on the left-hand side of the condition, and vice versa. To support this, column names can be prefixed by x\$ or y\$ to explicitly specify which table they come from.

```
sales <- tibble(</pre>
  id = c(1L, 1L, 1L, 2L, 2L),
 sale_date = as.Date(c("2018-12-31", "2019-01-02", "2019-01-05", "2019-01-04", "2019-01-01"))
)
sales
promos <- tibble(</pre>
  id = c(1L, 1L, 2L),
  promo_date = as.Date(c("2019-01-01", "2019-01-05", "2019-01-02"))
)
promos
# Match `id` to `id`, and `sale_date` to `promo_date`
by <- join_by(id, sale_date == promo_date)</pre>
left_join(sales, promos, by)
# For each `sale_date` within a particular `id`,
# find all `promo_date`s that occurred before that particular sale
by <- join_by(id, sale_date >= promo_date)
left_join(sales, promos, by)
# For each `sale_date` within a particular `id`,
# find only the closest `promo_date` that occurred before that sale
by <- join_by(id, closest(sale_date >= promo_date))
left_join(sales, promos, by)
# If you want to disallow exact matching in rolling joins, use `>` rather
# than `>=`. Note that the promo on 2019-01-05 is no longer considered the
# closest match for the sale on the same date.
by <- join_by(id, closest(sale_date > promo_date))
left_join(sales, promos, by)
```

```
# Same as before, but also require that the promo had to occur at most 1
# day before the sale was made. We'll use a full join to see that id 2's
# promo on 2019-01-02 is no longer matched to the sale on 2019-01-04.
sales <- mutate(sales, sale_date_lower = sale_date - 1)</pre>
by <- join_by(id, closest(sale_date >= promo_date), sale_date_lower <= promo_date)</pre>
full_join(sales, promos, by)
segments <- tibble(</pre>
 segment_id = 1:4,
 chromosome = c("chr1", "chr2", "chr2", "chr1"),
 start = c(140, 210, 380, 230),
 end = c(150, 240, 415, 280)
)
segments
reference <- tibble(</pre>
 reference_id = 1:4,
 chromosome = c("chr1", "chr1", "chr2", "chr2"),
 start = c(100, 200, 300, 415),
 end = c(150, 250, 399, 450)
)
reference
# Find every time a segment `start` falls between the reference
# `[start, end]` range.
by <- join_by(chromosome, between(start, start, end))</pre>
full_join(segments, reference, by)
# If you wanted the reference columns first, supply `reference` as `x`
# and `segments` as `y`, then explicitly refer to their columns using `x$`
# and `y$`.
by <- join_by(chromosome, between(y$start, x$start, x$end))</pre>
full_join(reference, segments, by)
# Find every time a segment falls completely within a reference.
# Sometimes using `x$` and `y$` makes your intentions clearer, even if they
# match the default behavior.
by <- join_by(chromosome, within(x$start, x$end, y$start, y$end))</pre>
inner_join(segments, reference, by)
# Find every time a segment overlaps a reference in any way.
by <- join_by(chromosome, overlaps(x$start, x$end, y$start, y$end))</pre>
full_join(segments, reference, by)
# It is common to have right-open ranges with bounds like `[)`, which would
# mean an end value of `415` would no longer overlap a start value of `415`.
# Setting `bounds` allows you to compute overlaps with those kinds of ranges.
by <- join_by(chromosome, overlaps(x$start, x$end, y$start, y$end, bounds = "[)"))</pre>
full_join(segments, reference, by)
```

lead-lag

## Description

Find the "previous" (lag()) or "next" (lead()) values in a vector. Useful for comparing values behind of or ahead of the current values.

### Usage

```
lag(x, n = 1L, default = NULL, order_by = NULL, ...)
lead(x, n = 1L, default = NULL, order_by = NULL, ...)
```

#### Arguments

x	A vector
n	Positive integer of length 1, giving the number of positions to lag or lead by
default	The value used to pad x back to its original size after the lag or lead has been applied. The default, NULL, pads with a missing value. If supplied, this must be a vector with size 1, which will be cast to the type of $x$ .
order_by	An optional secondary vector that defines the ordering to use when applying the lag or lead to x. If supplied, this must be the same size as x.
	Not used.

### Value

A vector with the same type and size as x.

```
lag(1:5)
lead(1:5)
x <- 1:5
tibble(behind = lag(x), x, ahead = lead(x))
# If you want to look more rows behind or ahead, use `n`
lag(1:5, n = 1)
lag(1:5, n = 2)
lead(1:5, n = 2)
# If you want to define a value to pad with, use `default`
lag(1:5)
lag(1:5, default = 0)
```

mutate

```
lead(1:5)
lead(1:5, default = 6)
# If the data are not already ordered, use `order_by`
scrambled <- slice_sample(
   tibble(year = 2000:2005, value = (0:5) ^ 2),
   prop = 1
)
wrong <- mutate(scrambled, previous_year_value = lag(value))
arrange(wrong, year)
right <- mutate(scrambled, previous_year_value = lag(value, order_by = year))
arrange(right, year)
```

```
mutate
```

Create, modify, and delete columns

#### Description

mutate() creates new columns that are functions of existing variables. It can also modify (if the name is the same as an existing column) and delete columns (by setting their value to NULL).

#### Usage

```
mutate(.data, ...)
## S3 method for class 'data.frame'
mutate(
   .data,
   ...,
   .by = NULL,
   .keep = c("all", "used", "unused", "none"),
   .before = NULL,
   .after = NULL
)
```

#### Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g.
	from dbplyr or dtplyr). See Methods, below, for more details.
	<pre><data-masking> Name-value pairs. The name gives the name of the column in</data-masking></pre>

<data-masking> Name-value pairs. The name gives the name of the column in
the output.

The value can be:

- A vector of length 1, which will be recycled to the correct length.
- A vector the same length as the current group (or the whole data frame if ungrouped).

	• NULL, to remove the column.
	• A data frame or tibble, to create multiple columns in the output.
.by	[Experimental]
	<tidy-select> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to group_by(). For details and examples, see ?dplyr_by.</tidy-select>
.keep	Control which columns from . data are retained in the output. Grouping columns and columns created by are always kept.
	• "all" retains all columns from . data. This is the default.
	• "used" retains only the columns used in to create new columns. This is useful for checking your work, as it displays inputs and outputs side-by-side.
	• "unused" retains only the columns <i>not</i> used in to create new columns. This is useful if you generate new columns, but no longer need the columns used to generate them.
	• "none" doesn't retain any extra columns from .data. Only the grouping variables and columns created by are kept.
.before, .after	<tidy-select> Optionally, control where new columns should appear (the default is to add to the right hand side). See relocate() for more details.</tidy-select>

#### Value

An object of the same type as .data. The output has the following properties:

- Columns from . data will be preserved according to the .keep argument.
- Existing columns that are modified by ... will always be returned in their original location.
- New columns created through . . . will be placed according to the . before and .after arguments.
- The number of rows is not affected.
- Columns given the value NULL will be removed.
- Groups will be recomputed if a grouping variable is mutated.
- Data frame attributes are preserved.

#### Useful mutate functions

- +, -, log(), etc., for their usual mathematical meanings
- lead(), lag()
- dense\_rank(), min\_rank(), percent\_rank(), row\_number(), cume\_dist(), ntile()
- cumsum(), cummean(), cummin(), cummax(), cumany(), cumall()
- na\_if(), coalesce()
- if\_else(), recode(), case\_when()

#### mutate

### Grouped tibbles

Because mutating expressions are computed within groups, they may yield different results on grouped tibbles. This will be the case as soon as an aggregating, lagging, or ranking function is involved. Compare this ungrouped mutate:

```
starwars %>%
  select(name, mass, species) %>%
  mutate(mass_norm = mass / mean(mass, na.rm = TRUE))
```

With the grouped equivalent:

```
starwars %>%
  select(name, mass, species) %>%
  group_by(species) %>%
  mutate(mass_norm = mass / mean(mass, na.rm = TRUE))
```

The former normalises mass by the global average whereas the latter normalises by the averages within species levels.

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages: no methods found.

#### See Also

Other single table verbs: arrange(), filter(), reframe(), rename(), select(), slice(), summarise()

```
# Newly created variables are available immediately
starwars %>%
  select(name, mass) %>%
  mutate(
    mass2 = mass * 2,
    mass2_squared = mass2 * mass2
)
# As well as adding new variables, you can use mutate() to
# remove variables and modify existing variables.
starwars %>%
  select(name, height, mass, homeworld) %>%
  mutate(
    mass = NULL,
    height = height * 0.0328084 # convert to feet
)
```

```
# Use across() with mutate() to apply a transformation
# to multiple columns in a tibble.
starwars %>%
 select(name, homeworld, species) %>%
 mutate(across(!name, as.factor))
# see more in ?across
# Window functions are useful for grouped mutates:
starwars %>%
 select(name, mass, homeworld) %>%
 group_by(homeworld) %>%
 mutate(rank = min_rank(desc(mass)))
# see `vignette("window-functions")` for more details
# By default, new columns are placed on the far right.
df <- tibble(x = 1, y = 2)
df %>% mutate(z = x + y)
df %>% mutate(z = x + y, .before = 1)
df %>% mutate(z = x + y, .after = x)
# By default, mutate() keeps all columns from the input data.
df <- tibble(x = 1, y = 2, a = "a", b = "b")
df %>% mutate(z = x + y, .keep = "all") # the default
df %>% mutate(z = x + y, .keep = "used")
df %>% mutate(z = x + y, .keep = "unused")
df %>% mutate(z = x + y, .keep = "none")
# Grouping ------
# The mutate operation may yield different results on grouped
# tibbles because the expressions are computed within groups.
# The following normalises `mass` by the global average:
starwars %>%
 select(name, mass, species) %>%
 mutate(mass_norm = mass / mean(mass, na.rm = TRUE))
# Whereas this normalises `mass` by the averages within species
# levels:
starwars %>%
 select(name, mass, species) %>%
 group_by(species) %>%
 mutate(mass_norm = mass / mean(mass, na.rm = TRUE))
# Indirection ------
# Refer to column names stored as strings with the `.data` pronoun:
vars <- c("mass", "height")</pre>
mutate(starwars, prod = .data[[vars[[1]]]] * .data[[vars[[2]]]])
# Learn more in ?rlang::args_data_masking
```

mutate-joins

Mutating joins

58

#### mutate-joins

### Description

Mutating joins add columns from y to x, matching observations based on the keys. There are four mutating joins: the inner join, and the three outer joins.

### Inner join:

An inner\_join() only keeps observations from x that have a matching key in y.

The most important property of an inner join is that unmatched rows in either input are not included in the result. This means that generally inner joins are not appropriate in most analyses, because it is too easy to lose observations.

### **Outer joins:**

The three outer joins keep observations that appear in at least one of the data frames:

- A left\_join() keeps all observations in x.
- A right\_join() keeps all observations in y.
- A full\_join() keeps all observations in x and y.

#### Usage

```
inner_join(
 х,
 у,
 by = NULL,
 copy = FALSE,
  suffix = c(".x", ".y"),
  . . . ,
 keep = NULL
)
## S3 method for class 'data.frame'
inner_join(
 х,
 у,
 by = NULL,
  copy = FALSE,
  suffix = c(".x", ".y"),
  . . . ,
 keep = NULL,
  na_matches = c("na", "never"),
 multiple = "all",
 unmatched = "drop"
  relationship = NULL
)
left_join(
 х,
 у,
 by = NULL,
 copy = FALSE,
```

```
suffix = c(".x", ".y"),
  ...,
 keep = NULL
)
## S3 method for class 'data.frame'
left_join(
 х,
 у,
 by = NULL,
 copy = FALSE,
  suffix = c(".x", ".y"),
  ...,
 keep = NULL,
  na_matches = c("na", "never"),
 multiple = "all",
 unmatched = "drop",
  relationship = NULL
)
right_join(
 х,
 у,
 by = NULL,
 copy = FALSE,
 suffix = c(".x", ".y"),
  ...,
 keep = NULL
)
## S3 method for class 'data.frame'
right_join(
 х,
 у,
 by = NULL,
  copy = FALSE,
 suffix = c(".x", ".y"),
  . . . ,
  keep = NULL,
  na_matches = c("na", "never"),
 multiple = "all",
 unmatched = "drop",
  relationship = NULL
)
full_join(
 х,
```

```
у,
```

60

## mutate-joins

```
by = NULL,
 copy = FALSE,
 suffix = c(".x", ".y"),
  · · · ,
 keep = NULL
)
## S3 method for class 'data.frame'
full_join(
 х,
 у,
 by = NULL,
 copy = FALSE,
 suffix = c(".x", ".y"),
  ...,
  keep = NULL,
 na_matches = c("na", "never"),
 multiple = "all",
 relationship = NULL
)
```

## Arguments

х, у	A pair of data frames, data frame extensions (e.g. a tibble), or lazy data frames (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
by	A join specification created with join_by(), or a character vector of variables to join by.
	If NULL, the default, *_join() will perform a natural join, using all variables in common across x and y. A message lists the variables so that you can check they're correct; suppress the message by supplying by explicitly.
	To join on different variables between x and y, use a join_by() specification. For example, join_by(a == b) will match x\$a to y\$b.
	To join by multiple variables, use a join_by() specification with multiple expressions. For example, join_by(a == b, c == d) will match x\$a to y\$b and x\$c to y\$d. If the column names are the same between x and y, you can shorten this by listing only the variable names, like join_by(a, c).
	join_by() can also be used to perform inequality, rolling, and overlap joins. See the documentation at ?join_by for details on these types of joins.
	For simple equality joins, you can alternatively specify a character vector of variable names to join by. For example, by = $c("a", "b")$ joins x\$a to y\$a and x\$b to y\$b. If variable names differ between x and y, use a named character vector like by = $c("x_a" = "y_a", "x_b" = "y_b")$ .
	To perform a cross-join, generating all combinations of x and y, see cross_join().
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.
suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2.

	Other parameters passed onto methods.
keep	Should the join keys from both x and y be preserved in the output?
	<ul> <li>If NULL, the default, joins on equality retain only the keys from x, while joins on inequality retain the keys from both inputs.</li> <li>If TRUE, all keys from both inputs are retained.</li> <li>If FALSE, only keys from x are retained. For right and full joins, the data in key columns corresponding to rows that only exist in y are merged into the key columns from x. Can't be used when joining on inequality conditions.</li> </ul>
na_matches	Should two NA or two NaN values match?
	<ul> <li>"na", the default, treats two NA or two NaN values as equal, like %in%, match(), and merge().</li> <li>"never" treats two NA or two NaN values as different, and will never match them together or to any other values. This is similar to joins for database sources and to base::merge(incomparables = NA).</li> </ul>
multiple	Handling of rows in x with multiple matches in y. For each row of x:
	• "all", the default, returns every match detected in y. This is the same behavior as SQL.
	<ul> <li>"any" returns one match detected in y, with no guarantees on which match will be returned. It is often faster than "first" and "last" if you just need to detect if there is at least one match.</li> <li>"first" returns the first match detected in y.</li> </ul>
	<ul> <li>"last" returns the last match detected in y.</li> </ul>
unmatched	How should unmatched keys that would result in dropped rows be handled?
unina ceneu	<ul> <li>"drop" drops unmatched keys from the result</li> </ul>
	<ul> <li>"error" throws an error if unmatched keys are detected.</li> </ul>
	unmatched is intended to protect you from accidentally dropping rows during a join. It only checks for unmatched keys in the input that could potentially drop rows.
	• For left joins, it checks y.
	• For right joins, it checks x.
	• For inner joins, it checks both x and y. In this case, unmatched is also allowed to be a character vector of length 2 to specify the behavior for x and y independently.
relationship	Handling of the expected relationship between the keys of x and y. If the expec- tations chosen from the list below are invalidated, an error is thrown.
	<ul> <li>NULL, the default, doesn't expect there to be any relationship between x and y. However, for equality joins it will check for a many-to-many relationship (which is typically unexpected) and will warn if one occurs, encouraging you to either take a closer look at your inputs or make this relationship explicit by specifying "many-to-many". See the <i>Many-to-many relationships</i> section for more details.</li> <li>"one-to-one" expects: <ul> <li>Each row in x matches at most 1 row in y.</li> </ul> </li> </ul>

- Each row in y matches at most 1 row in x.
- "one-to-many" expects:
  - Each row in y matches at most 1 row in x.
- "many-to-one" expects:
  - Each row in x matches at most 1 row in y.
- "many-to-many" doesn't perform any relationship checks, but is provided to allow you to be explicit about this relationship if you know it exists.

relationship doesn't handle cases where there are zero matches. For that, see unmatched.

#### Value

An object of the same type as x (including the same groups). The order of the rows and columns of x is preserved as much as possible. The output has the following properties:

- The rows are affect by the join type.
  - inner\_join() returns matched x rows.
  - left\_join() returns all x rows.
  - right\_join() returns matched of x rows, followed by unmatched y rows.
  - full\_join() returns all x rows, followed by unmatched y rows.
- Output columns include all columns from x and all non-key columns from y. If keep = TRUE, the key columns from y are included as well.
- If non-key columns in x and y have the same name, suffixes are added to disambiguate. If keep = TRUE and key columns in x and y have the same name, suffixes are added to disambiguate these as well.
- If keep = FALSE, output columns included in by are coerced to their common type between x and y.

### Many-to-many relationships

By default, dplyr guards against many-to-many relationships in equality joins by throwing a warning. These occur when both of the following are true:

- A row in x matches multiple rows in y.
- A row in y matches multiple rows in x.

This is typically surprising, as most joins involve a relationship of one-to-one, one-to-many, or many-to-one, and is often the result of an improperly specified join. Many-to-many relationships are particularly problematic because they can result in a Cartesian explosion of the number of rows returned from the join.

If a many-to-many relationship is expected, silence this warning by explicitly setting relationship = "many-to-many".

In production code, it is best to preemptively set relationship to whatever relationship you expect to exist between the keys of x and y, as this forces an error to occur immediately if the data doesn't align with your expectations.

Inequality joins typically result in many-to-many relationships by nature, so they don't warn on them by default, but you should still take extra care when specifying an inequality join, because they also have the capability to return a large number of rows.

Rolling joins don't warn on many-to-many relationships either, but many rolling joins follow a many-to-one relationship, so it is often useful to set relationship = "many-to-one" to enforce this.

Note that in SQL, most database providers won't let you specify a many-to-many relationship between two tables, instead requiring that you create a third *junction table* that results in two one-tomany relationships instead.

#### Methods

These functions are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

- inner\_join(): no methods found.
- left\_join(): no methods found.
- right\_join(): no methods found.
- full\_join(): no methods found.

#### See Also

Other joins: cross\_join(), filter-joins, nest\_join()

```
band_members %>% inner_join(band_instruments)
band_members %>% left_join(band_instruments)
band_members %>% right_join(band_instruments)
band_members %>% full_join(band_instruments)
# To suppress the message about joining variables, supply `by`
band_members %>% inner_join(band_instruments, by = join_by(name))
# This is good practice in production code
# Use an equality expression if the join variables have different names
band_members %>% full_join(band_instruments2, by = join_by(name == artist))
# By default, the join keys from `x` and `y` are coalesced in the output; use
# `keep = TRUE` to keep the join keys from both `x` and `y`
band_members %>%
  full_join(band_instruments2, by = join_by(name == artist), keep = TRUE)
# If a row in `x` matches multiple rows in `y`, all the rows in `y` will be
# returned once for each matching row in `x`.
df1 <- tibble(x = 1:3)
df2 <- tibble(x = c(1, 1, 2), y = c("first", "second", "third"))
df1 %>% left_join(df2)
```

```
# If a row in `y` also matches multiple rows in `x`, this is known as a
# many-to-many relationship, which is typically a result of an improperly
# specified join or some kind of messy data. In this case, a warning is
# thrown by default:
df3 <- tibble(x = c(1, 1, 1, 3))
df3 %>% left_join(df2)
# In the rare case where a many-to-many relationship is expected, set
# `relationship = "many-to-many"` to silence this warning
df3 %>% left_join(df2, relationship = "many-to-many")
# Use `join_by()` with a condition other than `==` to perform an inequality
# join. Here we match on every instance where df1x > df2x.
df1 %>% left_join(df2, join_by(x > x))
# By default, NAs match other NAs so that there are two
# rows in the output of this join:
df1 <- data.frame(x = c(1, NA), y = 2)
df2 <- data.frame(x = c(1, NA), z = 3)
left_join(df1, df2)
# You can optionally request that NAs don't match, giving a
# a result that more closely resembles SQL joins
left_join(df1, df2, na_matches = "never")
```

na if	Convert values to NA
110_11	

### Description

This is a translation of the SQL command NULLIF. It is useful if you want to convert an annoying value to NA.

#### Usage

na\_if(x, y)

#### Arguments

х	Vector to modify
У	Value or vector to compare against. When $x$ and $y$ are equal, the value in $x$ will be replaced with NA.
	y is cast to the type of x before comparison.
	y is recycled to the size of x before comparison. This means that y can be a vector with the same size as $x$ , but most of the time this will be a single value.

#### Value

A modified version of x that replaces any values that are equal to y with NA.

## See Also

coalesce() to replace missing values with a specified value.

tidyr::replace\_na() to replace NA with a value.

# Examples

```
na_if(1:5, 5:1)
x <- c(1, -1, 0, 10)
100 / x
100 / na_if(x, 0)
y <- c("abc", "def", "", "ghi")</pre>
na_if(y, "")
# `na_if()` allows you to replace `NaN` with `NA`,
# even though `NaN == NaN` returns `NA`
z <- c(1, NaN, NA, 2, NaN)
na_if(z, NaN)
# `na_if()` is particularly useful inside `mutate()`,
# and is meant for use with vectors rather than entire data frames
starwars %>%
  select(name, eye_color) %>%
  mutate(eye_color = na_if(eye_color, "unknown"))
# `na_if()` can also be used with `mutate()` and `across()`
# to alter multiple columns
starwars %>%
   mutate(across(where(is.character), ~na_if(., "unknown")))
```

near

Compare two numeric vectors

## Description

This is a safe way of comparing if two vectors of floating point numbers are (pairwise) equal. This is safer than using ==, because it has a built in tolerance

### Usage

near(x, y, tol = .Machine\$double.eps^0.5)

### Arguments

х, у	Numeric vectors to compare
tol	Tolerance of comparison.

near

## nest\_join

## Examples

sqrt(2) ^ 2 == 2
near(sqrt(2) ^ 2, 2)

nest\_join

Nest join

## Description

A nest join leaves x almost unchanged, except that it adds a new list-column, where each element contains the rows from y that match the corresponding row in x.

## Usage

```
nest_join(x, y, by = NULL, copy = FALSE, keep = NULL, name = NULL, ...)
## S3 method for class 'data.frame'
nest_join(
    x,
    y,
    by = NULL,
    copy = FALSE,
    keep = NULL,
    name = NULL,
    ...,
    na_matches = c("na", "never"),
    unmatched = "drop"
)
```

### Arguments

х, у	A pair of data frames, data frame extensions (e.g. a tibble), or lazy data frames (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
by	A join specification created with join_by(), or a character vector of variables to join by.
	If NULL, the default, *_join() will perform a natural join, using all variables in common across x and y. A message lists the variables so that you can check they're correct; suppress the message by supplying by explicitly.
	To join on different variables between x and y, use a join_by() specification. For example, join_by(a == b) will match x\$a to y\$b.
	To join by multiple variables, use a join_by() specification with multiple expressions. For example, join_by(a == b, c == d) will match x\$a to y\$b and x\$c to y\$d. If the column names are the same between x and y, you can shorten this by listing only the variable names, like join_by(a, c).
	join_by() can also be used to perform inequality, rolling, and overlap joins. See the documentation at ?join_by for details on these types of joins.

	For simple equality joins, you can alternatively specify a character vector of variable names to join by. For example, by = $c("a", "b")$ joins x\$a to y\$a and x\$b to y\$b. If variable names differ between x and y, use a named character vector like by = $c("x_a" = "y_a", "x_b" = "y_b")$ .
	To perform a cross-join, generating all combinations of x and y, see cross_join().
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.
keep	Should the new list-column contain join keys? The default will preserve the join keys for inequality joins.
name	The name of the list-column created by the join. If NULL, the default, the name of y is used.
	Other parameters passed onto methods.
na_matches	Should two NA or two NaN values match?
	<ul> <li>"na", the default, treats two NA or two NaN values as equal, like %in%, match(), and merge().</li> </ul>
	<ul> <li>"never" treats two NA or two NaN values as different, and will never match them together or to any other values. This is similar to joins for database sources and to base::merge(incomparables = NA).</li> </ul>
unmatched	How should unmatched keys that would result in dropped rows be handled?
	<ul> <li>"drop" drops unmatched keys from the result.</li> <li>"error" throws an error if unmatched keys are detected</li> </ul>
	unmatched is intended to protect you from accidentally dropping rows during a join. It only checks for unmatched keys in the input that could potentially drop rows.
	• For left joins, it checks y.
	• For right joins, it checks x.
	• For inner joins, it checks both x and y. In this case, unmatched is also allowed to be a character vector of length 2 to specify the behavior for x and y independently.

# Value

The output:

- Is same type as x (including having the same groups).
- Has exactly the same number of rows as x.
- Contains all the columns of x in the same order with the same values. They are only modified (slightly) if keep = FALSE, when columns listed in by will be coerced to their common type across x and y.
- Gains one new column called {name} on the far right, a list column containing data frames the same type as y.

### **Relationship to other joins**

You can recreate many other joins from the result of a nest join:

- inner\_join() is a nest\_join() plus tidyr::unnest().
- left\_join() is a nest\_join() plus tidyr::unnest(keep\_empty = TRUE).
- semi\_join() is a nest\_join() plus a filter() where you check that every element of data
  has at least one row.
- anti\_join() is a nest\_join() plus a filter() where you check that every element has zero rows.

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

### See Also

Other joins: cross\_join(), filter-joins, mutate-joins

#### Examples

```
df1 <- tibble(x = 1:3)
df2 <- tibble(x = c(2, 3, 3), y = c("a", "b", "c"))
out <- nest_join(df1, df2)
out
out$
df2</pre>
```

nth

Extract the first, last, or nth value from a vector

#### Description

These are useful helpers for extracting a single value from a vector. They are guaranteed to return a meaningful value, even when the input is shorter than expected. You can also provide an optional secondary vector that defines the ordering.

## Usage

nth(x, n, order\_by = NULL, default = NULL, na\_rm = FALSE)
first(x, order\_by = NULL, default = NULL, na\_rm = FALSE)
last(x, order\_by = NULL, default = NULL, na\_rm = FALSE)

#### Arguments

х	A vector
n	For nth(), a single integer specifying the position. Negative integers index from the end (i.e. $-1L$ will return the last value in the vector).
order_by	An optional vector the same size as x used to determine the order.
default	A default value to use if the position does not exist in x. If NULL, the default, a missing value is used. If supplied, this must be a single value, which will be cast to the type of x. When x is a list , default is allowed to be any value. There are no type or size restrictions in this case.
na_rm	Should missing values in x be removed before extracting the value?

## Details

For most vector types, first(x), last(x), and nth(x, n) work like x[[1]], x[[length(x)], and x[[n]], respectively. The primary exception is data frames, where they instead retrieve rows, i.e. x[1, ], x[nrow(x), ], and x[n, ]. This is consistent with the tidyverse/vctrs principle which treats data frames as a vector of rows, rather than a vector of columns.

#### Value

If x is a list, a single element from that list. Otherwise, a vector the same type as x with size 1.

```
x <- 1:10
y <- 10:1
first(x)
last(y)
nth(x, 1)
nth(x, 5)
nth(x, -2)
# `first()` and `last()` are often useful in `summarise()`
df <- tibble(x = x, y = y)
df %>%
  summarise(
   across(x:y, first, .names = "{col}_first"),
   y_{last} = last(y)
  )
# Selecting a position that is out of bounds returns a default value
nth(x, 11)
nth(x, 0)
# This out of bounds behavior also applies to empty vectors
first(integer())
```

ntile

```
# You can customize the default value with `default`
nth(x, 11, default = -1L)
first(integer(), default = 0L)
# `order_by` provides optional ordering
last(x)
last(x, order_by = y)
# `na_rm` removes missing values before extracting the value
z <- c(NA, NA, 1, 3, NA, 5, NA)
first(z)
first(z, na_rm = TRUE)
last(z, na_rm = TRUE)
nth(z, 3, na_rm = TRUE)
# For data frames, these select entire rows
df <- tibble(a = 1:5, b = 6:10)
first(df)
nth(df, 4)
```

ntile

Bucket a numeric vector into n groups

#### Description

ntile() is a sort of very rough rank, which breaks the input vector into n buckets. If length(x) is not an integer multiple of n, the size of the buckets will differ by up to one, with larger buckets coming first.

Unlike other ranking functions, ntile() ignores ties: it will create evenly sized buckets even if the same value of x ends up in different buckets.

## Usage

ntile(x = row\_number(), n)

## Arguments

х	A vector to rank
	By default, the smallest values will get the smallest ranks. Use desc() to reverse the direction so the largest values get the smallest ranks.
	Missing values will be given rank NA. Use coalesce(x, Inf) or coalesce(x, -Inf) if you want to treat them as the largest or smallest values respectively.
	To rank by multiple columns at once, supply a data frame.
n	Number of groups to bucket into

## See Also

Other ranking functions: percent\_rank(), row\_number()

## n\_distinct

### Examples

```
x <- c(5, 1, 3, 2, 2, NA)
ntile(x, 2)
ntile(x, 4)
# If the bucket sizes are uneven, the larger buckets come first
ntile(1:8, 3)
# Ties are ignored
ntile(rep(1, 8), 3)</pre>
```

n\_distinct Count unique combinations

## Description

n\_distinct() counts the number of unique/distinct combinations in a set of one or more vectors. It's a faster and more concise equivalent to nrow(unique(data.frame(...))).

### Usage

n\_distinct(..., na.rm = FALSE)

#### Arguments

	Unnamed vectors. If multiple vectors are supplied, then they should have the same length.
na.rm	If TRUE, exclude missing observations from the count. If there are multiple vec-
	tors in, an observation will be excluded if <i>any</i> of the values are missing.

### Value

A single number.

## Examples

```
x <- c(1, 1, 2, 2, 2)
n_distinct(x)
y <- c(3, 3, NA, 3, 3)
n_distinct(y)
n_distinct(y, na.rm = TRUE)
# Pairs (1, 3), (2, 3), and (2, NA) are distinct
n_distinct(x, y)
# (2, NA) is dropped, leaving 2 distinct combinations
n_distinct(x, y, na.rm = TRUE)
```

72
## order\_by

```
# Also works with data frames
n_distinct(data.frame(x, y))
```

order\_by

A helper function for ordering window function output

# Description

This function makes it possible to control the ordering of window functions in R that don't have a specific ordering parameter. When translated to SQL it will modify the order clause of the OVER function.

# Usage

```
order_by(order_by, call)
```

# Arguments

order_by	a vector to order_by
call	a function call to a window function, where the first argument is the vector being
	operated on

# Details

This function works by changing the call to instead call with\_order() with the appropriate arguments.

### Examples

```
order_by(10:1, cumsum(1:10))
x <- 10:1
y <- 1:10
order_by(x, cumsum(y))

df <- data.frame(year = 2000:2005, value = (0:5) ^ 2)
scrambled <- df[sample(nrow(df)), ]

wrong <- mutate(scrambled, running = cumsum(value))
arrange(wrong, year)

right <- mutate(scrambled, running = order_by(year, cumsum(value)))
arrange(right, year)</pre>
```

percent\_rank

# Description

These two ranking functions implement two slightly different ways to compute a percentile. For each x\_i in x:

- cume\_dist(x) counts the total number of values less than or equal to x\_i, and divides it by the number of observations.
- percent\_rank(x) counts the total number of values less than x\_i, and divides it by the number of observations minus 1.

In both cases, missing values are ignored when counting the number of observations.

### Usage

percent\_rank(x)

cume\_dist(x)

## Arguments

```
х
```

A vector to rank

By default, the smallest values will get the smallest ranks. Use desc() to reverse the direction so the largest values get the smallest ranks. Missing values will be given rank NA. Use coalesce(x, Inf) or coalesce(x, -Inf) if you want to treat them as the largest or smallest values respectively. To rank by multiple columns at once, supply a data frame.

#### Value

A numeric vector containing a proportion.

### See Also

Other ranking functions: ntile(), row\_number()

### Examples

```
x <- c(5, 1, 3, 2, 2)
cume_dist(x)
percent_rank(x)
# You can understand what's going on by computing it by hand
sapply(x, function(xi) sum(x <= xi) / length(x))</pre>
```

sapply(x, function(xi) sum(x < xi) / (length(x) - 1))

```
# The real computations are a little more complex in order to
# correctly deal with missing values
```

pick

## Select a subset of columns

## Description

pick() provides a way to easily select a subset of columns from your data using select() semantics while inside a "data-masking" function like mutate() or summarise(). pick() returns a data frame containing the selected columns for the current group.

pick() is complementary to across():

- With pick(), you typically apply a function to the full data frame.
- With across(), you typically apply a function to each column.

#### Usage

pick(...)

#### Arguments

• • •

# <tidy-select>

Columns to pick.

You can't pick grouping columns because they are already automatically handled by the verb (i.e. summarise() or mutate()).

## Details

Theoretically, pick() is intended to be replaceable with an equivalent call to tibble(). For example, pick(a, c) could be replaced with tibble(a = a, c = c), and pick(everything()) on a data frame with cols a, b, and c could be replaced with tibble(a = a, b = b, c = c). pick() specially handles the case of an empty selection by returning a 1 row, 0 column tibble, so an exact replacement is more like:

```
size <- vctrs::vec_size_common(..., .absent = 1L)
out <- vctrs::vec_recycle_common(..., .size = size)
tibble::new_tibble(out, nrow = size)</pre>
```

# Value

A tibble containing the selected columns for the current group.

## See Also

across()

# Examples

```
df <- tibble(</pre>
  x = c(3, 2, 2, 2, 1),
  y = c(0, 2, 1, 1, 4),
 z1 = c("a", "a", "a", "b", "a"),
  z2 = c("c", "d", "d", "a", "c")
)
df
# `pick()` provides a way to select a subset of your columns using
# tidyselect. It returns a data frame.
df %>% mutate(cols = pick(x, y))
# This is useful for functions that take data frames as inputs.
# For example, you can compute a joint rank between `x` and `y`.
df %>% mutate(rank = dense_rank(pick(x, y)))
# `pick()` is also useful as a bridge between data-masking functions (like
# `mutate()` or `group_by()`) and functions with tidy-select behavior (like
# `select()`). For example, you can use `pick()` to create a wrapper around
# `group_by()` that takes a tidy-selection of columns to group on. For more
# bridge patterns, see
# https://rlang.r-lib.org/reference/topic-data-mask-programming.html#bridge-patterns.
my_group_by <- function(data, cols) {</pre>
  group_by(data, pick({{ cols }}))
}
df %>% my_group_by(c(x, starts_with("z")))
# Or you can use it to dynamically select columns to `count()` by
df %>% count(pick(starts_with("z")))
```

pull

#### Extract a single column

# Description

pull() is similar to \$. It's mostly useful because it looks a little nicer in pipes, it also works with remote data frames, and it can optionally name the output.

#### Usage

pull(.data, var = -1, name = NULL, ...)

### Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g
	from dbplyr or dtplyr). See Methods, below, for more details.
var	A variable specified as:

76

# recode

	• a literal variable name
	• a positive integer, giving the position counting from the left
	• a negative integer, giving the position counting from the right.
	The default returns the last column (on the assumption that's the column you've created most recently).
	This argument is taken by expression and supports quasiquotation (you can un- quote column names and column locations).
name	An optional parameter that specifies the column to be used as names for a named vector. Specified in a similar manner as var.
	For use by methods.

# Value

A vector the same size as .data.

# Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

# Examples

```
mtcars %>% pull(-1)
mtcars %>% pull(1)
mtcars %>% pull(cyl)

# Also works for remote sources
df <- dbplyr::memdb_frame(x = 1:10, y = 10:1, .name = "pull-ex")
df %>%
  mutate(z = x * y) %>%
  pull()

# Pull a named vector
starwars %>% pull(height, name)
```



Recode values

# Description

## [Superseded]

recode() is superseded in favor of case\_match(), which handles the most important cases of recode() with a more elegant interface. recode\_factor() is also superseded, however, its direct replacement is not currently available but will eventually live in forcats. For creating new variables based on logical vectors, use if\_else(). For even more complicated criteria, use case\_when().

recode() is a vectorised version of switch(): you can replace numeric values based on their position or their name, and character or factor values only by their name. This is an S3 generic: dplyr provides methods for numeric, character, and factors. You can use recode() directly with factors; it will preserve the existing order of levels while changing the values. Alternatively, you can use recode\_factor(), which will change the order of levels to match the order of replacements.

#### Usage

```
recode(.x, ..., .default = NULL, .missing = NULL)
recode_factor(.x, ..., .default = NULL, .missing = NULL, .ordered = FALSE)
```

#### Arguments

. X	A vector to modify
	<pre><dynamic-dots> Replacements. For character and factor .x, these should be named and replacement is based only on their name. For numeric .x, these can be named or not. If not named, the replacement is done based on position i.ex represents positions to look for in replacements. See examples.</dynamic-dots></pre>
	When named, the argument names should be the current values to be replaced, and the argument values should be the new (replacement) values.
	All replacements must be the same type, and must have either length one or the same length as .x.
.default	If supplied, all values not otherwise matched will be given this value. If not supplied and if the replacements are the same type as the original values in .x, unmatched values are not changed. If not supplied and if the replacements are not compatible, unmatched values are replaced with NA.
	.default must be either length 1 or the same length as .x.
.missing	If supplied, any missing values in $.x$ will be replaced by this value. Must be either length 1 or the same length as $.x$ .
.ordered	If TRUE, recode_factor() creates an ordered factor.

# Value

A vector the same length as .x, and the same type as the first of ..., .default, or .missing. recode\_factor() returns a factor whose levels are in the same order as in .... The levels in .default and .missing come last.

#### recode

## See Also

na\_if() to replace specified values with a NA.

coalesce() to replace missing values with a specified value.

```
tidyr::replace_na() to replace NA with a value.
```

## Examples

```
char_vec <- sample(c("a", "b", "c"), 10, replace = TRUE)</pre>
# `recode()` is superseded by `case_match()`
recode(char_vec, a = "Apple", b = "Banana")
case_match(char_vec, "a" ~ "Apple", "b" ~ "Banana", .default = char_vec)
# With `case_match()`, you don't need typed missings like `NA_character_`
recode(char_vec, a = "Apple", b = "Banana", .default = NA_character_)
case_match(char_vec, "a" ~ "Apple", "b" ~ "Banana", .default = NA)
# Throws an error as `NA` is logical, not character.
try(recode(char_vec, a = "Apple", b = "Banana", .default = NA))
# case_match() is easier to use with numeric vectors, because you don't
# need to turn the numeric values into names
num_vec <- c(1:4, NA)
recode(num_vec, `2` = 20L, `4` = 40L)
case_match(num_vec, 2 \sim 20, 4 \sim 40, .default = num_vec)
# `case_match()` doesn't have the ability to match by position like
# recode() does with numeric vectors
recode(num_vec, "a", "b", "c", "d")
recode(c(1,5,3), "a", "b", "c", "d", .default = "nothing")
# For `case_match()`, incompatible types are an error rather than a warning
recode(num_vec, `2` = "b", `4` = "d")
try(case_match(num_vec, 2 ~ "b", 4 ~ "d", .default = num_vec))
# The factor method of `recode()` can generally be replaced with
# `forcats::fct_recode()`
factor_vec <- factor(c("a", "b", "c"))</pre>
recode(factor_vec, a = "Apple")
# `recode_factor()` does not currently have a direct replacement, but we
# plan to add one to forcats. In the meantime, you can use the `.ptype`
# argument to `case_match()`.
recode_factor(
  num_vec,
  `1` = "z"
  `2` = "y",
  `3` = "x",
  .default = "D"
  .missing = "M"
)
```

# reframe

```
case_match(
    num_vec,
    1 ~ "z",
    2 ~ "y",
    3 ~ "x",
    NA ~ "M",
    .default = "D",
    .ptype = factor(levels = c("z", "y", "x", "D", "M"))
)
```

reframe

Transform each group to an arbitrary number of rows

# Description

## [Experimental]

While summarise() requires that each argument returns a single value, and mutate() requires that each argument returns the same number of rows as the input, reframe() is a more general workhorse with no requirements on the number of rows returned per group.

reframe() creates a new data frame by applying functions to columns of an existing data frame. It is most similar to summarise(), with two big differences:

- reframe() can return an arbitrary number of rows per group, while summarise() reduces each group down to a single row.
- reframe() always returns an ungrouped data frame, while summarise() might return a grouped or rowwise data frame, depending on the scenario.

We expect that you'll use summarise() much more often than reframe(), but reframe() can be particularly helpful when you need to apply a complex function that doesn't return a single summary value.

#### Usage

```
reframe(.data, ..., .by = NULL)
```

## Arguments

A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
<data-masking></data-masking>
Name-value pairs of functions. The name will be the name of the variable in the result. The value can be a vector of any length.
Unnamed data frame values add multiple columns from a single expression.
[Experimental]
<tidy-select> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to group_by(). For details and examples, see ?dplyr_by.</tidy-select>

80

# reframe

# Value

If . data is a tibble, a tibble. Otherwise, a data.frame.

- The rows originate from the underlying grouping keys.
- The columns are a combination of the grouping keys and the expressions that you provide.
- The output is always ungrouped.
- Data frame attributes are **not** preserved, because reframe() fundamentally creates a new data frame.

# **Connection to tibble**

reframe() is theoretically connected to two functions in tibble: :enframe() and tibble::deframe():

- enframe(): vector -> data frame
- deframe(): data frame -> vector
- reframe(): data frame -> data frame

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

## See Also

Other single table verbs: arrange(), filter(), mutate(), rename(), select(), slice(), summarise()

# Examples

```
table <- c("a", "b", "d", "f")

df <- tibble(
  g = c(1, 1, 1, 2, 2, 2, 2),
  x = c("e", "a", "b", "c", "f", "d", "a")
)

# `reframe()` allows you to apply functions that return
# an arbitrary number of rows
df %>%
  reframe(x = intersect(x, table))

# Functions are applied per group, and each group can return a
# different number of rows.
df %>%
  reframe(x = intersect(x, table), .by = g)
# The output is always ungrouped, even when using `group_by()`
df %>%
```

```
relocate
```

```
group_by(g) %>%
  reframe(x = intersect(x, table))
# You can add multiple columns at once using a single expression by returning
# a data frame.
quantile_df <- function(x, probs = c(0.25, 0.5, 0.75)) {</pre>
  tibble(
   val = quantile(x, probs, na.rm = TRUE),
    quant = probs
 )
}
x <- c(10, 15, 18, 12)
quantile_df(x)
starwars %>%
  reframe(quantile_df(height))
starwars %>%
  reframe(quantile_df(height), .by = homeworld)
starwars %>%
  reframe(
   across(c(height, mass), quantile_df, .unpack = TRUE),
    .by = homeworld
  )
```

relocate

#### Change column order

# Description

Use relocate() to change column positions, using the same syntax as select() to make it easy to move blocks of columns at once.

# Usage

```
relocate(.data, ..., .before = NULL, .after = NULL)
```

# Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<tidy-select> Columns to move.</tidy-select>
.before, .after	<tidy-select> Destination of columns selected by Supplying neither will move columns to the left-hand side; specifying both is an error.</tidy-select>

82

#### rename

## Value

An object of the same type as . data. The output has the following properties:

- Rows are not affected.
- The same columns appear in the output, but (usually) in a different place and possibly renamed.
- Data frame attributes are preserved.
- Groups are not affected.

## Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

# Examples

```
df <- tibble(a = 1, b = 1, c = 1, d = "a", e = "a", f = "a")
df %>% relocate(f)
df %>% relocate(a, .after = c)
df %>% relocate(f, .before = b)
df %>% relocate(a, .after = last_col())
# relocated columns can change name
df %>% relocate(ff = f)
# Can also select variables based on their type
df %>% relocate(where(is.character))
df %>% relocate(where(is.numeric), .after = last_col())
# Or with any other select helper
df %>% relocate(any_of(c("a", "e", "i", "o", "u")))
# When .before or .after refers to multiple variables they will be
# moved to be immediately before/after the selected variables.
df2 <- tibble(a = 1, b = "a", c = 1, d = "a")
df2 %>% relocate(where(is.numeric), .after = where(is.character))
```

```
df2 %>% relocate(where(is.numeric), .before = where(is.character))
```

rename

Rename columns

#### Description

rename() changes the names of individual variables using new\_name = old\_name syntax; rename\_with()
renames columns using a function.

rename

### Usage

```
rename(.data, ...)
```

rename\_with(.data, .fn, .cols = everything(), ...)

#### Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	For rename(): <tidy-select> Use new_name = old_name to rename selected variables.</tidy-select>
	For rename_with(): additional arguments passed onto .fn.
.fn	A function used to transform the selected .cols. Should return a character vector the same length as the input.
.cols	<tidy-select> Columns to rename; defaults to all columns.</tidy-select>

# Value

An object of the same type as .data. The output has the following properties:

- Rows are not affected.
- Column names are changed; column order is preserved.
- Data frame attributes are preserved.
- Groups are updated to reflect new names.

# Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

# See Also

Other single table verbs: arrange(), filter(), mutate(), reframe(), select(), slice(), summarise()

# Examples

```
iris <- as_tibble(iris) # so it prints a little nicer
rename(iris, petal_length = Petal.Length)
# Rename using a named vector and `all_of()`
lookup <- c(pl = "Petal.Length", sl = "Sepal.Length")
rename(iris, all_of(lookup))
# If your named vector might contain names that don't exist in the data,
# use `any_of()` instead
lookup <- c(lookup, new = "unknown")</pre>
```

84

rows

```
try(rename(iris, all_of(lookup)))
rename(iris, any_of(lookup))
rename_with(iris, toupper)
rename_with(iris, toupper, starts_with("Petal"))
rename_with(iris, ~ tolower(gsub(".", "_", .x, fixed = TRUE)))
# If your renaming function uses `paste0()`, make sure to set
# `recycle0 = TRUE` to ensure that empty selections are recycled correctly
try(rename_with(
 iris,
 ~ paste0("prefix_", .x),
 starts_with("nonexistent")
))
rename_with(
 iris,
 ~ paste0("prefix_", .x, recycle0 = TRUE),
 starts_with("nonexistent")
)
```

rows

#### Manipulate individual rows

## Description

These functions provide a framework for modifying rows in a table using a second table of data. The two tables are matched by a set of key variables whose values typically uniquely identify each row. The functions are inspired by SQL's INSERT, UPDATE, and DELETE, and can optionally modify in\_place for selected backends.

- rows\_insert() adds new rows (like INSERT). By default, key values in y must not exist in x.
- rows\_append() works like rows\_insert() but ignores keys.
- rows\_update() modifies existing rows (like UPDATE). Key values in y must be unique, and, by default, key values in y must exist in x.
- rows\_patch() works like rows\_update() but only overwrites NA values.
- rows\_upsert() inserts or updates depending on whether or not the key value in y already exists in x. Key values in y must be unique.
- rows\_delete() deletes rows (like DELETE). By default, key values in y must exist in x.

### Usage

```
rows_insert(
    x,
    y,
```

```
by = NULL,
  ...,
  conflict = c("error", "ignore"),
  copy = FALSE,
 in_place = FALSE
)
rows_append(x, y, ..., copy = FALSE, in_place = FALSE)
rows_update(
 х,
 у,
 by = NULL,
  ...,
 unmatched = c("error", "ignore"),
  copy = FALSE,
  in_place = FALSE
)
rows_patch(
 х,
 у,
 by = NULL,
  ...,
 unmatched = c("error", "ignore"),
 copy = FALSE,
  in_place = FALSE
)
rows_upsert(x, y, by = NULL, ..., copy = FALSE, in_place = FALSE)
rows_delete(
 х,
 у,
 by = NULL,
  ...,
 unmatched = c("error", "ignore"),
 copy = FALSE,
  in_place = FALSE
)
```

#### Arguments

х, у	A pair of data frames or data frame extensions (e.g. a tibble). y must have the
	same columns of x or a subset.

by An unnamed character vector giving the key columns. The key columns must exist in both x and y. Keys typically uniquely identify each row, but this is only enforced for the key values of y when rows\_update(), rows\_patch(), or

86

	rows_upsert() are used.
	By default, we use the first column in y, since the first column is a reasonable place to put an identifier variable.
	Other parameters passed onto methods.
conflict	For rows_insert(), how should keys in y that conflict with keys in x be han- dled? A conflict arises if there is a key in y that already exists in x.
	One of:
	• "error", the default, will error if there are any keys in y that conflict with keys in x.
	• "ignore" will ignore rows in y with keys that conflict with keys in x.
сору	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. This allows you to join tables across srcs, but it is a potentially expensive operation so you must opt into it.
in_place	Should x be modified in place? This argument is only relevant for mutable backends (e.g. databases, data.tables).
	When TRUE, a modified version of x is returned invisibly; when FALSE, a new object representing the resulting changes is returned.
unmatched	For rows_update(), rows_patch(), and rows_delete(), how should keys in y that are unmatched by the keys in x be handled?
	One of:
	• "error", the default, will error if there are any keys in y that are unmatched by the keys in x.
	• "ignore" will ignore rows in y with keys that are unmatched by the keys in x.

#### Value

An object of the same type as x. The order of the rows and columns of x is preserved as much as possible. The output has the following properties:

- rows\_update() and rows\_patch() preserve the number of rows; rows\_insert(), rows\_append(), and rows\_upsert() return all existing rows and potentially new rows; rows\_delete() returns a subset of the rows.
- Columns are not added, removed, or relocated, though the data may be updated.
- Groups are taken from x.
- Data frame attributes are taken from x.

If in\_place = TRUE, the result will be returned invisibly.

### Methods

These function are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

- rows\_insert(): no methods found.
- rows\_append(): no methods found.
- rows\_update(): no methods found.
- rows\_patch(): no methods found.
- rows\_upsert(): no methods found.
- rows\_delete(): no methods found.

## Examples

```
data <- tibble(a = 1:3, b = letters[c(1:2, NA)], c = 0.5 + 0:2)</pre>
data
# Insert
rows_insert(data, tibble(a = 4, b = "z"))
# By default, if a key in `y` matches a key in `x`, then it can't be inserted
# and will throw an error. Alternatively, you can ignore rows in `y`
# containing keys that conflict with keys in `x` with `conflict = "ignore"`,
# or you can use `rows_append()` to ignore keys entirely.
try(rows_insert(data, tibble(a = 3, b = "z")))
rows_insert(data, tibble(a = 3, b = "z"), conflict = "ignore")
rows_append(data, tibble(a = 3, b = "z"))
# Update
rows_update(data, tibble(a = 2:3, b = "z"))
rows_update(data, tibble(b = "z", a = 2:3), by = "a")
# Variants: patch and upsert
rows_patch(data, tibble(a = 2:3, b = "z"))
rows_upsert(data, tibble(a = 2:4, b = "z"))
# Delete and truncate
rows_delete(data, tibble(a = 2:3))
rows_delete(data, tibble(a = 2:3, b = "b"))
# By default, for update, patch, and delete it is an error if a key in `y`
# doesn't exist in `x`. You can ignore rows in `y` that have unmatched keys
# with `unmatched = "ignore"`.
y <- tibble(a = 3:4, b = "z")</pre>
try(rows_update(data, y, by = "a"))
rows_update(data, y, by = "a", unmatched = "ignore")
rows_patch(data, y, by = "a", unmatched = "ignore")
rows_delete(data, y, by = "a", unmatched = "ignore")
```

rowwise

#### rowwise

# Description

rowwise() allows you to compute on a data frame a row-at-a-time. This is most useful when a vectorised function doesn't exist.

Most dplyr verbs preserve row-wise grouping. The exception is summarise(), which return a grouped\_df. You can explicitly ungroup with ungroup() or as\_tibble(), or convert to a grouped\_df with group\_by().

### Usage

rowwise(data, ...)

## Arguments

data	Input data frame.
	<tidy-select> Variable</tidy-select>
	typically a set of variable

es to be preserved when calling summarise(). This is s whose combination uniquely identify each row. NB: unlike group\_by() you can not create new variables here but instead you can select multiple variables with (e.g.) everything().

## Value

A row-wise data frame with class rowwise\_df. Note that a rowwise\_df is implicitly grouped by row, but is not a grouped\_df.

## List-columns

Because a rowwise has exactly one row per group it offers a small convenience for working with list-columns. Normally, summarise() and mutate() extract a groups worth of data with [. But when you index a list in this way, you get back another list. When you're working with a rowwise tibble, then dplyr will use [[ instead of [ to make your life a little easier.

### See Also

nest\_by() for a convenient way of creating rowwise data frames with nested data.

## Examples

```
df <- tibble(x = runif(6), y = runif(6), z = runif(6))
# Compute the mean of x, y, z in each row
df %>% rowwise() %>% mutate(m = mean(c(x, y, z)))
# use c_across() to more easily select many variables
df %>% rowwise() %>% mutate(m = mean(c_across(x:z)))
# Compute the minimum of x and y in each row
df %>% rowwise() %>% mutate(m = min(c(x, y, z)))
# In this case you can use an existing vectorised function:
df %>% mutate(m = pmin(x, y, z))
# Where these functions exist they'll be much faster than rowwise
# so be on the lookout for them.
```

```
# rowwise() is also useful when doing simulations
params <- tribble(</pre>
 ~sim, ~n, ~mean, ~sd,
   1, 1,
             1, 1,
   2, 2,
              2,
                    4,
   3, 3,
             -1,
                    2
)
# Here I supply variables to preserve after the computation
params %>%
 rowwise(sim) %>%
 reframe(z = rnorm(n, mean, sd))
# If you want one row per simulation, put the results in a list()
params %>%
 rowwise(sim) %>%
 summarise(z = list(rnorm(n, mean, sd)), .groups = "keep")
```

```
row_number
```

Integer ranking functions

#### Description

Three ranking functions inspired by SQL2003. They differ primarily in how they handle ties:

- row\_number() gives every input a unique rank, so that c(10, 20, 20, 30) would get ranks c(1, 2, 3, 4). It's equivalent to rank(ties.method = "first").
- min\_rank() gives every tie the same (smallest) value so that c(10, 20, 20, 30) gets ranks c(1, 2, 2, 4). It's the way that ranks are usually computed in sports and is equivalent to rank(ties.method = "min").
- dense\_rank() works like min\_rank(), but doesn't leave any gaps, so that c(10, 20, 20, 30) gets ranks c(1, 2, 2, 3).

## Usage

```
row_number(x)
```

min\_rank(x)

dense\_rank(x)

#### Arguments

Х

A vector to rank

By default, the smallest values will get the smallest ranks. Use desc() to reverse the direction so the largest values get the smallest ranks. Missing values will be given rank NA. Use coalesce(x, Inf) or coalesce(x,

-Inf) if you want to treat them as the largest or smallest values respectively.

To rank by multiple columns at once, supply a data frame.

```
90
```

scoped

# Value

An integer vector.

## See Also

Other ranking functions: ntile(), percent\_rank()

# Examples

```
x <- c(5, 1, 3, 2, 2, NA)
row_number(x)
min_rank(x)
dense_rank(x)
# Ranking functions can be used in `filter()` to select top/bottom rows
df <- data.frame(</pre>
  grp = c(1, 1, 1, 2, 2, 2, 3, 3, 3),
  x = c(3, 2, 1, 1, 2, 2, 1, 1, 1),
  y = c(1, 3, 2, 3, 2, 2, 4, 1, 2),
  id = 1:9
)
# Always gives exactly 1 row per group
df %>% group_by(grp) %>% filter(row_number(x) == 1)
# May give more than 1 row if ties
df %>% group_by(grp) %>% filter(min_rank(x) == 1)
# Rank by multiple columns (to break ties) by selecting them with `pick()`
df %>% group_by(grp) %>% filter(min_rank(pick(x, y)) == 1)
# See slice_min() and slice_max() for another way to tackle the same problem
# You can use row_number() without an argument to refer to the "current"
# row number.
df %>% group_by(grp) %>% filter(row_number() == 1)
# It's easiest to see what this does with mutate():
df %>% group_by(grp) %>% mutate(grp_id = row_number())
```

scoped

Operate on a selection of variables

## Description

# [Superseded]

Scoped verbs (\_if, \_at, \_all) have been superseded by the use of pick() or across() in an existing verb. See vignette("colwise") for details.

The variants suffixed with \_if, \_at or \_all apply an expression (sometimes several) to all variables within a specified subset. This subset can contain all variables (\_all variants), a vars() selection (\_at variants), or variables selected with a predicate (\_if variants).

The verbs with scoped variants are:

- mutate(), transmute() and summarise(). See summarise\_all().
- filter(). See filter\_all().
- group\_by(). See group\_by\_all().
- rename() and select(). See select\_all().
- arrange(). See arrange\_all()

There are three kinds of scoped variants. They differ in the scope of the variable selection on which operations are applied:

- Verbs suffixed with \_all() apply an operation on all variables.
- Verbs suffixed with \_at() apply an operation on a subset of variables specified with the quoting function vars(). This quoting function accepts tidyselect::vars\_select() helpers like starts\_with(). Instead of a vars() selection, you can also supply an integerish vector of column positions or a character vector of column names.
- Verbs suffixed with \_if() apply an operation on the subset of variables for which a predicate function returns TRUE. Instead of a predicate function, you can also supply a logical vector.

#### Arguments

.tbl	A tbl object.
.funs	A function fun, a quosure style lambda ~ fun(.) or a list of either form.
.vars	A list of columns generated by vars(), a character vector of column names, a numeric vector of column positions, or NULL.
.predicate	A predicate function to be applied to the columns or a logical vector. The variables for which .predicate is or returns TRUE are selected. This argument is passed to rlang::as_function() and thus supports quosure-style lambda functions and strings representing function names.
	Additional arguments for the function calls in .funs. These are evaluated only once, with tidy dots support.

# **Grouping variables**

Most of these operations also apply on the grouping variables when they are part of the selection. This includes:

- arrange\_all(), arrange\_at(), and arrange\_if()
- distinct\_all(), distinct\_at(), and distinct\_if()
- filter\_all(), filter\_at(), and filter\_if()
- group\_by\_all(), group\_by\_at(), and group\_by\_if()
- select\_all(), select\_at(), and select\_if()

This is not the case for summarising and mutating variants where operations are *not* applied on grouping variables. The behaviour depends on whether the selection is **implicit** (all and if selections) or **explicit** (at selections). Grouping variables covered by explicit selections (with summarise\_at(), mutate\_at(), and transmute\_at()) are always an error. For implicit selections, the grouping variables are always ignored. In this case, the level of verbosity depends on the kind of operation:

### select

- Summarising operations (summarise\_all() and summarise\_if()) ignore grouping variables silently because it is obvious that operations are not applied on grouping variables.
- On the other hand it isn't as obvious in the case of mutating operations (mutate\_all(), mutate\_if(), transmute\_all(), and transmute\_if()). For this reason, they issue a message indicating which grouping variables are ignored.

select

#### Keep or drop columns using their names and types

## Description

Select (and optionally rename) variables in a data frame, using a concise mini-language that makes it easy to refer to variables based on their name (e.g. a:f selects all columns from a on the left to f on the right) or type (e.g. where(is.numeric) selects all numeric columns).

#### **Overview of selection features:**

Tidyverse selections implement a dialect of R where operators make it easy to select variables:

- : for selecting a range of consecutive variables.
- ! for taking the complement of a set of variables.
- & and | for selecting the intersection or the union of two sets of variables.
- c() for combining selections.

In addition, you can use selection helpers. Some helpers select specific columns:

- everything(): Matches all variables.
- last\_col(): Select last variable, possibly with an offset.
- group\_cols(): Select all grouping columns.

Other helpers select variables by matching patterns in their names:

- starts\_with(): Starts with a prefix.
- ends\_with(): Ends with a suffix.
- contains(): Contains a literal string.
- matches(): Matches a regular expression.
- num\_range(): Matches a numerical range like x01, x02, x03.

Or from variables stored in a character vector:

- all\_of(): Matches variable names in a character vector. All names must be present, otherwise an out-of-bounds error is thrown.
- any\_of(): Same as all\_of(), except that no error is thrown for names that don't exist.

Or using a predicate function:

• where(): Applies a function to all variables and selects those for which the function returns TRUE.

### Usage

select(.data, ...)

# Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
• • •	<tidy-select> One or more unquoted expressions separated by commas. Vari- able names can be used as if they were positions in the data frame, so expressions</tidy-select>
	like x : y can be used to select a range of variables.

# Value

An object of the same type as .data. The output has the following properties:

- Rows are not affected.
- Output columns are a subset of input columns, potentially with a different order. Columns will be renamed if new\_name = old\_name form is used.
- Data frame attributes are preserved.
- Groups are maintained; you can't select off grouping variables.

# Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

### Examples

Here we show the usage for the basic selection operators. See the specific help pages to learn about helpers like starts\_with().

The selection language can be used in functions like dplyr::select() or tidyr::pivot\_longer(). Let's first attach the tidyverse:

library(tidyverse)

```
# For better printing
iris <- as_tibble(iris)</pre>
```

Select variables by name:

```
starwars %>% select(height)
#> # A tibble: 87 x 1
#> height
#> <int>
#> 1 172
#> 2 167
#> 3 96
#> 4 202
#> # i 83 more rows
```

```
iris %>% pivot_longer(Sepal.Length)
#> # A tibble: 150 x 6
     Sepal.Width Petal.Length Petal.Width Species name
#>
                                                                 value
#>
           <dbl>
                        <dbl>
                                     <dbl> <fct>
                                                   <chr>
                                                                 <dbl>
#> 1
             3.5
                          1.4
                                       0.2 setosa Sepal.Length
                                                                   5.1
#> 2
             3
                          1.4
                                       0.2 setosa Sepal.Length
                                                                   4.9
             3.2
#> 3
                          1.3
                                       0.2 setosa Sepal.Length
                                                                   4.7
             3.1
                          1.5
                                       0.2 setosa Sepal.Length
#> 4
                                                                   4.6
#> # i 146 more rows
```

Select multiple variables by separating them with commas. Note how the order of columns is determined by the order of inputs:

```
starwars %>% select(homeworld, height, mass)
#> # A tibble: 87 x 3
#>
     homeworld height mass
#>
     <chr>
                <int> <dbl>
#> 1 Tatooine
                   172
                          77
#> 2 Tatooine
                   167
                          75
#> 3 Naboo
                   96
                          32
                   202
#> 4 Tatooine
                         136
#> # i 83 more rows
```

Functions like tidyr::pivot\_longer() don't take variables with dots. In this case use c() to select multiple variables:

```
iris %>% pivot_longer(c(Sepal.Length, Petal.Length))
#> # A tibble: 300 x 5
#>
     Sepal.Width Petal.Width Species name
                                                   value
#>
           <dbl>
                       <dbl> <fct>
                                     <chr>
                                                   <dbl>
#> 1
             3.5
                         0.2 setosa Sepal.Length
                                                     5.1
#> 2
             3.5
                         0.2 setosa Petal.Length
                                                     1.4
#> 3
             3
                                                     4.9
                         0.2 setosa Sepal.Length
#> 4
             3
                         0.2 setosa Petal.Length
                                                     1.4
#> # i 296 more rows
```

#### **Operators::**

The : operator selects a range of consecutive variables:

<pre>starwars %&gt;% select(name:mass)</pre>				
#>	#	A tibble: 87 $\boldsymbol{x}$	3	
#>		name	height	mass
#>		<chr></chr>	<int></int>	<dbl></dbl>
#>	1	Luke Skywalker	172	77
#>	2	C-3P0	167	75
#>	3	R2-D2	96	32
#>	4	Darth Vader	202	136
#>	#	i 83 more rows		

```
The ! operator negates a selection:
```

```
starwars %>% select(!(name:mass))
#> # A tibble: 87 x 11
#> hair_color skin_color eye_color birth_year sex gender
                                                               homeworld species
#>
    <chr>
               <chr>
                           <chr>
                                          <dbl> <chr> <chr>
                                                                 <chr>
                                                                           <chr>
#> 1 blond
                fair
                           blue
                                           19
                                                male masculine Tatooine Human
#> 2 <NA>
                                                none masculine Tatooine
               gold
                           yellow
                                          112
                                                                           Droid
                                                                           Droid
#> 3 <NA>
               white, blue red
                                           33
                                                none masculine Naboo
#> 4 none
               white
                           yellow
                                           41.9 male masculine Tatooine Human
#> # i 83 more rows
#> # i 3 more variables: films <list>, vehicles <list>, starships <list>
iris %>% select(!c(Sepal.Length, Petal.Length))
#> # A tibble: 150 x 3
#>
     Sepal.Width Petal.Width Species
#>
           <dbl>
                       <dbl> <fct>
#> 1
             3.5
                         0.2 setosa
#> 2
             3
                         0.2 setosa
#> 3
             3.2
                         0.2 setosa
#> 4
             3.1
                         0.2 setosa
#> # i 146 more rows
iris %>% select(!ends_with("Width"))
#> # A tibble: 150 x 3
#>
     Sepal.Length Petal.Length Species
#>
                         <dbl> <fct>
            <dbl>
#> 1
              5.1
                           1.4 setosa
#> 2
              4.9
                           1.4 setosa
#> 3
              4.7
                           1.3 setosa
#> 4
              4.6
                           1.5 setosa
#> # i 146 more rows
& and | take the intersection or the union of two selections:
iris %>% select(starts_with("Petal") & ends_with("Width"))
#> # A tibble: 150 x 1
    Petal.Width
#>
#>
           <dbl>
#> 1
             0.2
#> 2
             0.2
#> 3
             0.2
#> 4
             0.2
#> # i 146 more rows
iris %>% select(starts_with("Petal") | ends_with("Width"))
#> # A tibble: 150 x 3
#>
    Petal.Length Petal.Width Sepal.Width
#>
            <dbl>
                        <dbl>
                                     <dbl>
#> 1
              1.4
                                       3.5
                           0.2
```

setops

#>	2	1.4	0.2	3
#>	3	1.3	0.2	3.2
#>	4	1.5	0.2	3.1
#>	# i	146 more rows		

To take the difference between two selections, combine the & and ! operators:

# See Also

Other single table verbs: arrange(), filter(), mutate(), reframe(), rename(), slice(), summarise()

setops

Set operations

# Description

Perform set operations using the rows of a data frame.

- intersect(x, y) finds all rows in both x and y.
- union(x, y) finds all rows in either x or y, excluding duplicates.
- union\_all(x, y) finds all rows in either x or y, including duplicates.
- setdiff(x, y) finds all rows in x that aren't in y.
- symdiff(x, y) computes the symmetric difference, i.e. all rows in x that aren't in y and all rows in y that aren't in x.
- setequal(x, y) returns TRUE if x and y contain the same rows (ignoring order).

Note that intersect(), union(), setdiff(), and symdiff() remove duplicates in x and y.

#### Usage

```
intersect(x, y, ...)
union(x, y, ...)
union_all(x, y, ...)
setdiff(x, y, ...)
```

```
setequal(x, y, ...)
```

symdiff(x, y, ...)

## Arguments

х, у	Pair of compatible data frames. A pair of data frames is compatible if they have the same column names (possibly in different orders) and compatible types.
	These dots are for future extensions and must be empty.

# **Base functions**

intersect(), union(), setdiff(), and setequal() override the base functions of the same name in order to make them generic. The existing behaviour for vectors is preserved by providing default methods that call the base functions.

# Examples

```
df1 <- tibble(x = 1:3)
df2 <- tibble(x = 3:5)
intersect(df1, df2)
union(df1, df2)
union_all(df1, df2)
setdiff(df1, df2)
setdiff(df2, df1)
symdiff(df1, df2)
setequal(df1, df2)
setequal(df1, df1[3:1, ])
# Note that the following functions remove pre-existing duplicates:
df1 <- tibble(x = c(1:3, 3, 3))
df2 <- tibble(x = c(3:5, 5))
intersect(df1, df2)
union(df1, df2)
setdiff(df1, df2)
symdiff(df1, df2)
```

```
slice
```

Subset rows using their positions

### Description

slice() lets you index rows by their (integer) locations. It allows you to select, remove, and duplicate rows. It is accompanied by a number of helpers for common use cases:

- slice\_head() and slice\_tail() select the first or last rows.
- slice\_sample() randomly selects rows.
- slice\_min() and slice\_max() select rows with the smallest or largest values of a variable.

If .data is a grouped\_df, the operation will be performed on each group, so that (e.g.) slice\_head(df, n = 5) will select the first five rows in each group.

# Usage

```
slice(.data, ..., .by = NULL, .preserve = FALSE)
slice_head(.data, ..., n, prop, by = NULL)
slice_tail(.data, ..., n, prop, by = NULL)
slice_min(
  .data,
 order_by,
  ...,
 n,
  prop,
 by = NULL,
 with_ties = TRUE,
 na_rm = FALSE
)
slice_max(
  .data,
 order_by,
  ...,
 n,
 prop,
 by = NULL,
 with_ties = TRUE,
 na_rm = FALSE
)
```

slice\_sample(.data, ..., n, prop, by = NULL, weight\_by = NULL, replace = FALSE)

# Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<pre>For slice(): <data-masking> Integer row values.</data-masking></pre>
	Provide either positive values to keep, or negative values to drop. The values provided must be either all positive or all negative. Indices beyond the number of rows in the input are silently ignored.
	For slice_*(), these arguments are passed on to methods.

.by,by	[Experimental]
	<tidy-select> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to group_by(). For details and examples, see ?dplyr_by.</tidy-select>
.preserve	Relevant when the .data input is grouped. If .preserve = FALSE (the default), the grouping structure is recalculated based on the resulting data, otherwise the grouping is kept as is.
n, prop	Provide either n, the number of rows, or prop, the proportion of rows to select. If neither are supplied, $n = 1$ will be used. If n is greater than the number of rows in the group (or prop > 1), the result will be silently truncated to the group size. prop will be rounded towards zero to generate an integer number of rows. A negative value of n or prop will be subtracted from the group size. For example, $n = -2$ with a group of 5 rows will select $5 - 2 = 3$ rows; prop $= -0.25$ with 8 rows will select $8 * (1 - 0.25) = 6$ rows.
order_by	<pre><data-masking> Variable or function of variables to order by. To order by multiple variables, wrap them in a data frame or tibble.</data-masking></pre>
with_ties	Should ties be kept together? The default, TRUE, may return more rows than you request. Use FALSE to ignore ties, and return the first n rows.
na_rm	Should missing values in order_by be removed from the result? If FALSE, NA values are sorted to the end (like in arrange()), so they will only be included if there are insufficient non-missing values to reach n/prop.
weight_by	<data-masking> Sampling weights. This must evaluate to a vector of non-negative numbers the same length as the input. Weights are automatically stan- dardised to sum to 1.</data-masking>
replace	Should sampling be performed with (TRUE) or without (FALSE, the default) replacement.

# Details

Slice does not work with relational databases because they have no intrinsic notion of row order. If you want to perform the equivalent operation, use filter() and row\_number().

# Value

An object of the same type as .data. The output has the following properties:

- Each row may appear 0, 1, or many times in the output.
- Columns are not modified.
- Groups are not modified.
- Data frame attributes are preserved.

# Methods

These function are **generics**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

- slice(): no methods found.
- slice\_head(): no methods found.
- slice\_tail(): no methods found.
- slice\_min(): no methods found.
- slice\_max(): no methods found.
- slice\_sample(): no methods found.

# See Also

```
Other single table verbs: arrange(), filter(), mutate(), reframe(), rename(), select(),
summarise()
```

# Examples

```
# Similar to head(mtcars, 1):
mtcars %>% slice(1L)
# Similar to tail(mtcars, 1):
mtcars %>% slice(n())
mtcars %>% slice(5:n())
# Rows can be dropped with negative indices:
slice(mtcars, -(1:4))
# First and last rows based on existing order
mtcars %>% slice_head(n = 5)
mtcars %>% slice_tail(n = 5)
# Rows with minimum and maximum values of a variable
mtcars %>% slice_min(mpg, n = 5)
mtcars %>% slice_max(mpg, n = 5)
# slice_min() and slice_max() may return more rows than requested
# in the presence of ties.
mtcars %>% slice_min(cyl, n = 1)
# Use with_ties = FALSE to return exactly n matches
mtcars %>% slice_min(cyl, n = 1, with_ties = FALSE)
# Or use additional variables to break the tie:
mtcars %>% slice_min(tibble(cyl, mpg), n = 1)
# slice_sample() allows you to random select with or without replacement
mtcars %>% slice_sample(n = 5)
mtcars %>% slice_sample(n = 5, replace = TRUE)
# you can optionally weight by a variable - this code weights by the
# physical weight of the cars, so heavy cars are more likely to get
# selected
mtcars %>% slice_sample(weight_by = wt, n = 5)
# Group wise operation ------
df <- tibble(</pre>
 group = rep(c("a", "b", "c"), c(1, 2, 4)),
 x = runif(7)
```

```
)
# All slice helpers operate per group, silently truncating to the group
# size, so the following code works without error
df %>% group_by(group) %>% slice_head(n = 2)
# When specifying the proportion of rows to include non-integer sizes
# are rounded down, so group a gets 0 rows
df %>% group_by(group) %>% slice_head(prop = 0.5)
# Filter equivalents -------
# slice() expressions can often be written to use `filter()` and
# `row_number()`, which can also be translated to SQL. For many databases,
# you'll need to supply an explicit variable to use to compute the row number.
filter(mtcars, row_number() == 1L)
filter(mtcars, tow_number() == n())
filter(mtcars, between(row_number(), 5, n()))
```

sql

SQL escaping.

#### Description

These functions are critical when writing functions that translate R functions to sql functions. Typically a conversion function should escape all its inputs and return an sql object.

#### Usage

sql(...)

## Arguments

• • •

Character vectors that will be combined into a single SQL expression.

starwars

Starwars characters

#### Description

The original data, from SWAPI, the Star Wars API, https://swapi.py4e.com/, has been revised to reflect additional research into gender and sex determinations of characters.

### Usage

starwars

102

#### storms

# Format

A tibble with 87 rows and 14 variables:

name Name of the character

height Height (cm)

mass Weight (kg)

hair\_color,skin\_color,eye\_color Hair, skin, and eye colors

**birth\_year** Year born (BBY = Before Battle of Yavin)

- **sex** The biological sex of the character, namely male, female, hermaphroditic, or none (as in the case for Droids).
- **gender** The gender role or gender identity of the character as determined by their personality or the way they were programmed (as in the case for Droids).

homeworld Name of homeworld

species Name of species

films List of films the character appeared in

vehicles List of vehicles the character has piloted

starships List of starships the character has piloted

Storm tracks data

#### Examples

starwars

storms

# Description

This dataset is the NOAA Atlantic hurricane database best track data, https://www.nhc.noaa. gov/data/#hurdat. The data includes the positions and attributes of storms from 1975-2022. Storms from 1979 onward are measured every six hours during the lifetime of the storm. Storms in earlier years have some missing data.

# Usage

storms

# Format

A tibble with 19,537 observations and 13 variables:

name Storm Name

year,month,day Date of report

hour Hour of report (in UTC)

#### summarise

lat, long Location of storm center

**status** Storm classification (Tropical Depression, Tropical Storm, or Hurricane) **category** Saffir-Simpson hurricane category calculated from wind speed.

- NA: Not a hurricane
- 1: 64+ knots
- 2: 83+ knots
- 3: 96+ knots
- 4: 113+ knots
- 5: 137+ knots

wind storm's maximum sustained wind speed (in knots)

pressure Air pressure at the storm's center (in millibars)

- **tropicalstorm\_force\_diameter** Diameter (in nautical miles) of the area experiencing tropical storm strength winds (34 knots or above). Only available starting in 2004.
- **hurricane\_force\_diameter** Diameter (in nautical miles) of the area experiencing hurricane strength winds (64 knots or above). Only available starting in 2004.

## See Also

The script to create the storms data set: https://github.com/tidyverse/dplyr/blob/main/ data-raw/storms.R

# Examples

storms

```
# Show a few recent storm paths
if (requireNamespace("ggplot2", quietly = TRUE)) {
    library(ggplot2)
    storms %>%
    filter(year >= 2000) %>%
    ggplot(aes(long, lat, color = paste(year, name))) +
    geom_path(show.legend = FALSE) +
    facet_wrap(~year)
}
storms
```

summarise

Summarise each group down to one row

#### Description

summarise() creates a new data frame. It returns one row for each combination of grouping variables; if there are no grouping variables, the output will have a single row summarising all observations in the input. It will contain one column for each grouping variable and one column for each of the summary statistics that you have specified.

summarise() and summarize() are synonyms.

# summarise

# Usage

<pre>summarise(.data,</pre>	· · · ,	.by = NULL,	.groups = NULL)
<pre>summarize(.data,</pre>	,	.by = NULL,	.groups = NULL)

# Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>Methods</i> , below, for more details.
	<data-masking> Name-value pairs of summary functions. The name will be the name of the variable in the result.</data-masking>
	The value can be:
	<ul> <li>A vector of length 1, e.g. min(x), n(), or sum(is.na(y)).</li> </ul>
	• A data frame, to add multiple columns from a single expression.
	[ <b>Deprecated</b> ] Returning values with size 0 or >1 was deprecated as of 1.1.0. Please use reframe() for this instead.
.by	[Experimental]
	<tidy-select> Optionally, a selection of columns to group by for just this operation, functioning as an alternative to group_by(). For details and examples, see ?dplyr_by.</tidy-select>
.groups	[Experimental] Grouping structure of the result.
	• "drop_last": dropping the last level of grouping. This was the only supported option before version 1.0.0.
	<ul> <li>"drop": All levels of grouping are dropped.</li> </ul>
	<ul> <li>"keep": Same grouping structure as .data.</li> </ul>
	• "rowwise": Each row is its own group.
	When .groups is not specified, it is chosen based on the number of rows of the results:
	• If all the results have 1 row, you get "drop_last".
	• If the number of rows varies, you get "keep" (note that returning a variable number of rows was deprecated in favor of reframe(), which also unconditionally drops all levels of grouping).
	In addition, a message informs you of that choice, unless the result is ungrouped, the option "dplyr.summarise.inform" is set to FALSE, or when summarise() is called from a function in a package.
Value	

An object usually of the same type as . data.

- The rows come from the underlying group\_keys().
- The columns are a combination of the grouping keys and the summary expressions that you provide.
- The grouping structure is controlled by the .groups= argument, the output may be another grouped\_df, a tibble or a rowwise data frame.

• Data frame attributes are **not** preserved, because summarise() fundamentally creates a new data frame.

# **Useful functions**

- Center: mean(), median()
- Spread: sd(), IQR(), mad()
- Range: min(), max(),
- Position: first(), last(), nth(),
- Count: n(), n\_distinct()
- Logical: any(), all()

## **Backend variations**

The data frame backend supports creating a variable and using it in the same summary. This means that previously created summary variables can be further transformed or combined within the summary, as in mutate(). However, it also means that summary variables with the same names as previous variables overwrite them, making those variables unavailable to later summary variables.

This behaviour may not be supported in other backends. To avoid unexpected results, consider using new names for your summary variables, especially when creating multiple summaries.

#### Methods

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages: no methods found.

### See Also

```
Other single table verbs: arrange(), filter(), mutate(), reframe(), rename(), select(),
slice()
```

# Examples

```
# A summary applied to ungrouped tbl returns a single row
mtcars %>%
   summarise(mean = mean(disp), n = n())
# Usually, you'll want to group first
mtcars %>%
   group_by(cyl) %>%
   summarise(mean = mean(disp), n = n())
# Each summary call removes one grouping level (since that group
# is now just a single row)
mtcars %>%
   group_by(cyl, vs) %>%
   summarise(cyl_n = n()) %>%
```

```
group_vars()
# BEWARE: reusing variables may lead to unexpected results
mtcars %>%
  group_by(cyl) %>%
  summarise(disp = mean(disp), sd = sd(disp))
# Refer to column names stored as strings with the `.data` pronoun:
var <- "mass"
summarise(starwars, avg = mean(.data[[var]], na.rm = TRUE))
# Learn more in ?rlang::args_data_masking
# In dplyr 1.1.0, returning multiple rows per group was deprecated in favor
# of `reframe()`, which never messages and always returns an ungrouped
# result:
mtcars %>%
   group_by(cyl) %>%
   summarise(qs = quantile(disp, c(0.25, 0.75)), prob = c(0.25, 0.75))
# ->
mtcars %>%
  group_by(cyl) %>%
   reframe(qs = quantile(disp, c(0.25, 0.75)), prob = c(0.25, 0.75))
```

tbl

#### *Create a table from a data source*

# Description

This is a generic method that dispatches based on the first argument.

#### Usage

tbl(src, ...)
is.tbl(x)

# Arguments

src	A data source
•••	Other arguments passed on to the individual methods
x	Any object

# Description

# [Superseded]

vars() is superseded because it is only needed for the scoped verbs (i.e. mutate\_at(), summarise\_at(), and friends), which have been been superseded in favour of across(). See vignette("colwise") for details.

This helper is intended to provide tidy-select semantics for scoped verbs like mutate\_at() and summarise\_at(). Note that anywhere you can supply vars() specification, you can also supply a numeric vector of column positions or a character vector of column names.

### Usage

vars(...)

# Arguments

... <tidy-select> Variables to operate on.

# See Also

all\_vars() and any\_vars() for other quoting functions that you can use with scoped verbs.

vars
# Index

\* datasets band\_members, 10 starwars, 102 storms, 103 \* grouping functions group\_by, 42 group\_map, 45 group\_trim, 47 \* joins cross\_join, 26 filter-joins, 39 mutate-joins, 58 nest\_join, 67 \* ranking functions ntile, 71 percent\_rank, 74 row\_number, 90 \* single table verbs arrange, 8 filter, 37 mutate, 55 reframe, 80 rename. 83 select, 93 slice, 98 summarise, 104 +, 56 ==, 38 >, 38 >=, 38 ?dplyr\_by, 37, 56, 80, 100, 105 ?join\_by, 40, 61, 67 &, <u>38</u> across, 3 across(), 7, 22, 29, 75, 91, 108 add\_count (count), 24 add\_tally (count), 24 all(), 106 all\_of(), 93

all\_vars,7 all\_vars(), 108 anti\_join (filter-joins), 39 anti\_join(), 69 any(), 106 any\_of(), 93 any\_vars(all\_vars), 7 any\_vars(), 108 arrange, 8, 38, 57, 81, 84, 97, 101, 106 arrange(), 29, 35, 43, 92, 100 arrange\_all(), 92 arrange\_at(), 92 arrange\_if(), 92 as\_tibble(), 89 auto\_copy, 10 band\_instruments (band\_members), 10 band\_instruments2 (band\_members), 10 band\_members, 10 between, 11 between(), 38bind (bind\_rows), 13 bind\_cols, 12 bind\_rows, 13 c\_across, 28 c\_across(), 5 case\_match, 14 case\_match(), 17, 78 case\_when, 16 case\_when(), 14, 15, 56, 78 cast, 65 closest (join\_by), 50 coalesce, 19 coalesce(), 56, 66, 79 collapse (compute), 20 collect (compute), 20 collect(), 24 compute, 20

consecutive\_id, 21

#### INDEX

contains(), 93 context, 22 copy\_to, 23 copy\_to(), 21 count, 24count(), 35cross\_join, 26, 41, 64, 69 cross\_join(), 40, 51, 61, 68 cumal1, 27 cumal1(), 56 cumany (cumall), 27 cumany(), 56cume\_dist (percent\_rank), 74 cume\_dist(), 56 cummax(), 56cummean (cumall), 27 cummean(), 56cummin(), 56cumsum(), 56cur\_column (context), 22 cur\_column(), 4 cur\_group (context), 22 cur\_group(), 4 cur\_group\_id (context), 22 cur\_group\_rows (context), 22 data-masking, 32 dense\_rank (row\_number), 90 dense\_rank(), 56 desc, 29 desc(), 8, 71, 74, 90 distinct. 30 distinct(), 35distinct\_all(), 92 distinct\_at(), 92 distinct\_if(), 92 do(), 46 dplyr-locale, 8 dplyr\_by, 31 ends\_with(), 93 everything(), 93 explain, 36 filter, 9, 37, 57, 81, 84, 97, 101, 106 filter(), 47, 92, 100 filter-joins, 39

filter\_all(), 7, 92
filter\_at(), 92

filter\_if(), 7, 92 first (nth), 69 first(), 106 full\_join (mutate-joins), 58 glimpse, 41 group\_by, 42, 46, 48 group\_by(), 25, 31, 35, 37, 48, 56, 80, 89, 92, 100.105 group\_by\_all(), 92 group\_by\_at(), 92 group\_by\_drop\_default(), 42 group\_by\_if(), 92 group\_cols, 44 group\_cols(), 93 group\_data(), 22 group\_keys(), 46, 105 group\_map, 43, 45, 48 group\_modify (group\_map), 45 group\_nest, 43, 46, 48 group\_split, 43, 46, 48 group\_trim, 43, 46, 47 group\_vars(), 45 group\_walk (group\_map), 45 grouped data frame, 48grouped\_df, 42, 89, 99, 105 groups(), 45ident, 48 if-else, 49 if\_all (across), 3 if\_any (across), 3 if\_else, 49 if\_else(), 16, 56, 78 ifelse(), 49

if\_all (across), 3
if\_any (across), 3
if\_else, 49
if\_else(), 16, 56, 78
ifelse(), 49
inner\_join (mutate-joins), 58
inner\_join(), 69
integerish, 92
intersect (setops), 97
IQR(), 106
is.na(), 38
is.tbl (tbl), 107

join, 12
join (mutate-joins), 58
join\_by, 50
join\_by(), 11, 40, 61, 67

lag (lead-lag), 54

## 110

## INDEX

lag(), 56 last (nth), 69 last(), 106 last\_col(), 93 lead (lead-lag), 54 lead(), 56 lead-lag, 54 left\_join (mutate-joins), 58 left\_join(), 50, 69 locale, 9 log(), 56 mad(), 106 match(), 40, 62, 68 matches(), 93max(), 106 mean(), 106 median(), 106 merge(), 40, 62, 68 min(), 106 min\_rank (row\_number), 90 min\_rank(), 56 mutate, 9, 38, 55, 81, 84, 97, 101, 106 mutate(), 3, 4, 22, 29, 75, 80, 92, 106 mutate-joins, 58 mutate\_all(), 93 mutate\_at(), 92, 108 mutate\_if(), 93 mutating joins, 26 n (context), 22 n(), 106 n\_distinct, 72 n\_distinct(), 106 na\_if, 65 na\_if(), 19, 56, 79 near, 66 near(), <u>38</u>  $nest_by(), 89$ nest\_join, 27, 41, 64, 67 nested, 46 nth, 69 nth(), 106 ntile, 71, 74, 91 ntile(), 56 num\_range(), 93 order\_by, 73 overlaps (join\_by), 50

percent\_rank, 71, 74, 91 percent\_rank(), 56 pick, 75 pick(), 3, 22, 23, 91 pillar::glimpse(), 41 print(), 36 pull, 76 pull(), 35 quasiquotation, 77 recode, 77 recode(), 56 recode\_factor (recode), 77 recycled, 12, 14, 19, 49, 65 reframe, 9, 38, 57, 80, 84, 97, 101, 106 reframe(), 105 relocate, 82 relocate(), 56 rename, 9, 38, 57, 81, 83, 97, 101, 106 rename(), 35, 92 rename\_with (rename), 83 right\_join (mutate-joins), 58 rlang::as\_function(), 92 row\_number, 71, 74, 90 row\_number(), 56, 100 rows.85 rows\_append (rows), 85 rows\_delete (rows), 85 rows\_insert (rows), 85 rows\_patch (rows), 85 rows\_update (rows), 85 rows\_upsert (rows), 85 rowwise, 88, 105 rowwise(), 29 scoped, 91 sd(), 106 select, 9, 38, 57, 81, 84, 93, 101, 106 select(), 3, 35, 44, 75, 92 select\_all(), 92 select\_at(), 92 select\_if(), 92 semi\_join(filter-joins), 39 semi\_join(), 69 setdiff (setops), 97 setequal (setops), 97 setops, 97

show\_query (explain), 36

### INDEX

slice, 9, 38, 57, 81, 84, 97, 98, 106 slice\_head (slice), 98 slice\_max (slice), 98 slice\_min(slice), 98 slice\_sample (slice), 98 slice\_tail (slice), 98 split, 46 sql, 102 starts\_with(), 92-94 starwars, 102 storms, 103str(), 36, 41 stringi::stri\_locale\_list(),9 summarise, 9, 38, 57, 81, 84, 97, 101, 104 summarise(), 3, 4, 22, 29, 32, 43, 75, 80, 89, 92 summarise\_all(), 92, 93 summarise\_at(), 92, 108 summarise\_if(), 93 summarize (summarise), 104 switch(), 14, 78 symdiff(setops), 97 tally (count), 24 tbl, 107 tbl(), 42 tibble, *105* tibble::deframe(), 81 tibble::enframe(), 81 tidy dots, 92 tidy-select, *32*, *34* tidyr::replace\_na(), 19, 66, 79 tidyr::unnest(), 69 tidyselect::vars\_select(), 92 transmute(), 92 transmute\_all(), 93 transmute\_at(), 92 transmute\_if(), 93 ungroup (group\_by), 42 ungroup(), 32, 37, 89 union (setops), 97 union\_all (setops), 97 unique.data.frame(), 30 unpack, 4 vars, 108 vars(), 8, 44, 91, 92

vctrs::vec\_as\_names(), 12

```
vctrs::vec_c(), 29
where(), 93
with_order(), 73
within(join_by), 50
```

xor(), <u>38</u>

112