## Package 'DFplyr'

October 24, 2025

```
Title A `DataFrame` (`S4Vectors`) backend for `dplyr`
Version 1.3.2
Description Provides 'dplyr' verbs ('mutate', 'select', 'filter', etc...)
      supporting `S4Vectors::DataFrame` objects. Importantly, this is achieved
      without conversion to an intermediate `tibble`. Adds grouping
     infrastructure to `DataFrame` which is respected by the transformation
      verbs.
biocViews DataRepresentation, Infrastructure, Software
License GPL-3
Encoding UTF-8
RoxygenNote 7.3.2
URL https://github.com/jonocarroll/DFplyr
BugReports https://github.com/jonocarroll/DFplyr/issues
Depends dplyr
Imports BiocGenerics, methods, rlang, S4Vectors, tidyselect
Suggests BiocStyle, GenomeInfoDb, GenomicRanges, IRanges, knitr,
     rmarkdown, sessioninfo, testthat (>= 3.0.0), tibble
VignetteBuilder knitr
Config/testthat/edition 3
Roxygen list(markdown = TRUE)
git_url https://git.bioconductor.org/packages/DFplyr
git_branch devel
git_last_commit 2aa58c3
git_last_commit_date 2025-09-26
Repository Bioconductor 3.22
Date/Publication 2025-10-24
Author Jonathan Carroll [aut, cre] (ORCID:
       <https://orcid.org/0000-0002-1404-5264>),
     Pierre-Paul Axisa [ctb]
Maintainer Jonathan Carroll <rpkg@jcarroll.com.au>
```

DFplyr-package

### **Contents**

	DFplyr-package	2
	arrange.DataFrame	3
	bindROWS,DataFrame-method	5
	count.DataFrame	5
	desc	7
	distinct.DataFrame	8
	filter.DataFrame	9
	format.DataFrame	11
	group_by.DataFrame	13
	group_by_drop_default.DataFrame	16
	group_data	17
	group_data.DataFrame	17
	group_vars.DataFrame	18
	inner_join.DataFrame	19
	mutate.DataFrame	20
	pull.DataFrame	22
	rename,DataFrame-method	23
	rename2	24
	select.DataFrame	24
	slice.DataFrame	28
	summarise.DataFrame	31
	summarize.DataFrame	33
	tally.DataFrame	35
	tbl vars.DataFrame	36
	<del>-</del>	37
	[,DataFrame-method	39
Index		41
DF 1	Total CAVactors DetaFrame of delice J	
n-bT	yr-package Treat a S4Vectors::DataFrame as a dplyr data source	

### Description

Add **dplyr** compatibility to S4Vectors::DataFrame for use with a selection of **dplyr** verbs.

### Arguments

x A S4Vectors::DataFrame object

### Author(s)

Maintainer: Jonathan Carroll <rpkg@jcarroll.com.au> (ORCID)

Other contributors:

• Pierre-Paul Axisa [contributor]

arrange.DataFrame 3

#### See Also

Useful links:

- https://github.com/jonocarroll/DFplyr
- Report bugs at https://github.com/jonocarroll/DFplyr/issues

### **Examples**

```
library(S4Vectors)
library(dplyr)
d <- as(mtcars, "DataFrame")</pre>
mutate(d, newvar = cyl + hp)
mutate_at(d, vars(starts_with("c")), ~ .^2)
group_by(d, cyl, am) %>%
    tally(gear)
count(d, gear, am, cyl)
select(d, am, cyl)
select(d, am, cyl) %>%
    rename2(foo = am)
arrange(d, desc(hp))
rbind(DataFrame(mtcars[1, ], row.names = "MyCar"), d) %>%
    distinct()
filter(d, am == 0)
slice(d, 3:6)
```

arrange.DataFrame

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

```
## S3 method for class 'DataFrame'
arrange(.data, ...)
```

4 arrange.DataFrame

#### **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>

#### **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))
```

bindROWS,DataFrame-method

rbind DataFrames

### Description

rbind DataFrames

### Usage

```
## S4 method for signature 'DataFrame'
bindROWS(x, objects = list())
```

### **Arguments**

```
x a DataFrameobjects a list of DataFrames
```

### Value

a new DataFrame combining the inputs by rows

count.DataFrame

Count the observations in each group

### **Description**

count() lets you quickly count the unique values of one or more variables: df %>% count(a, b) is roughly equivalent to df %>% group\_by(a, b) %>% summarise(n = n()). count() is paired with tally(), a lower-level helper that is equivalent to df %>% 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

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

6 count.DataFrame

#### **Arguments**

A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr).

... <data-masking> Variables to group by.

wt <data-masking> Frequency weights. Can be NULL or a variable:

- 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).

[Deprecated] 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",
# 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)
```

desc 7

```
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)
```

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**

Χ

vector to transform

#### Value

the input vector in a format that will be sorted in descending order.

```
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))
```

8 distinct.DataFrame

distinct.DataFrame Keep d

*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

```
## S3 method for class 'DataFrame'
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.
•••	<data-masking> Optional variables to use when determining uniqueness. If 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.</data-masking>
.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.

```
df <- tibble(
    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)</pre>
```

filter.DataFrame 9

```
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()
```

filter.DataFrame

Keep rows that match a condition

#### **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

```
## S3 method for class 'DataFrame'
filter(.data, ..., .preserve = FALSE)
```

### Arguments

8-----

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.

. . .

.data

<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 & operator. Only rows for which all conditions evaluate to TRUE are kept.

.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.

10 filter.DataFrame

#### **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.

format.DataFrame 11

#### 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(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]],
    . \verb|data[[vars[[2]]]]| > \verb|cond[[2]]||
# Learn more in ?rlang::args_data_masking
```

format.DataFrame

Encode in a Common Format

#### **Description**

Format an R object for pretty printing.

### Usage

```
## S3 method for class 'DataFrame' format(x, \ldots)
```

### Arguments

```
x any R object (conceptually); typically numeric.
```

... further arguments passed to or from other methods.

12 format.DataFrame

#### **Details**

format is a generic function. Apart from the methods described here there are methods for dates (see format.Date), date-times (see format.POSIXct) and for other classes such as format.octmode and format.dist.

format.data.frame formats the data frame column by column, applying the appropriate method of format for each column. Methods for columns are often similar to as.character but offer more control. Matrix and data-frame columns will be converted to separate columns in the result, and character columns (normally all) will be given class "AsIs".

format.factor converts the factor to a character vector and then calls the default method (and so justify applies).

format. AsIs deals with columns of complicated objects that have been extracted from a data frame. Character objects and (atomic) matrices are passed to the default method (and so width does not apply). Otherwise it calls toString to convert the object to character (if a vector or list, element by element) and then right-justifies the result.

Justification for character vectors (and objects converted to character vectors by their methods) is done on display width (see nchar), taking double-width characters and the rendering of special characters (as escape sequences, including escaping backslash but not double quote: see print.default) into account. Thus the width is as displayed by print(quote = FALSE) and not as displayed by cat. Character strings are padded with blanks to the display width of the widest. (If na.encode = FALSE missing character strings are not included in the width computations and are not encoded.)

Numeric vectors are encoded with the minimum number of decimal places needed to display all the elements to at least the digits significant digits. However, if all the elements then have trailing zeroes, the number of decimal places is reduced until at least one element has a non-zero final digit; see also the argument documentation for big.\*, small.\* etc, above. See the note in print.default about digits >= 16.

Raw vectors are converted to their 2-digit hexadecimal representation by as.character.

format.default(x) now provides a "minimal" string when isS4(x) is true.

While the internal code respects the option getOption("OutDec") for the 'decimal mark' in general, decimal.mark takes precedence over that option. Similarly, scientific takes precedence over getOption("scipen").

#### Value

An object of similar structure to x containing character representations of the elements of the first argument x in a common format, and in the current locale's encoding.

For character, numeric, complex or factor x, dims and dimnames are preserved on matrices/arrays and names on vectors: no other attributes are copied.

If x is a list, the result is a character vector obtained by applying format.default(x, ...) to each element of the list (after unlisting elements which are themselves lists), and then collapsing the result for each element with paste(collapse = ", "). The defaults in this case are trim = TRUE, justify = "none" since one does not usually want alignment in the collapsed strings.

### References

Becker, R. A., Chambers, J. M. and Wilks, A. R. (1988) *The New S Language*. Wadsworth & Brooks/Cole.

group\_by.DataFrame 13

#### See Also

format.info indicates how an atomic vector would be formatted. formatC, paste, as.character, sprintf, print, prettyNum, toString, encodeString.

### **Examples**

```
format(1:10)
format(1:10, trim = TRUE)
zz <- data.frame("(row names)"= c("aaaaa", "b"), check.names = FALSE)</pre>
format(zz)
format(zz, justify = "left")
## use of nsmall
format(13.7)
format(13.7, nsmall = 3)
format(c(6.0, 13.1), digits = 2)
format(c(6.0, 13.1), digits = 2, nsmall = 1)
## use of scientific
format(2<sup>31-1</sup>)
format(2^31-1, scientific = TRUE)
## scientific = numeric scipen (= {sci}entific notation {pen}alty) :
x <- c(1e5, 1000, 10, 0.1, .001, .123)
t(sapply(setNames(,-4:1),
         \(sci) sapply(x, format, scientific=sci)))
## a list
z \leftarrow list(a = letters[1:3], b = (-pi+0i)^((-2:2)/2), c = c(1,10,100,1000),
          d = c("a", "longer", "character", "string"),
          q = quote(a + b), e = expression(1+x))
## can you find the "2" small differences?
(f1 <- format(z, digits = 2))</pre>
(f2 <- format(z, digits = 2, justify = "left", trim = FALSE))</pre>
f1 == f2 ## 2 FALSE, 4 TRUE
## A "minimal" format() for S4 objects without their own format() method:
cc <- methods::getClassDef("standardGeneric")</pre>
format(cc) ## "<S4 class .....>"
```

group\_by.DataFrame

Group by one or more variables

### 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

```
## S3 method for class 'DataFrame'
group_by(.data, ..., add = FALSE, .drop = group_by_drop_default(.data))
```

group\_by.DataFrame

#### **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.

#### 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.

#### **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

group\_by.DataFrame 15

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()
```

```
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
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(
    x = 1:10,
    y = factor(rep(c("a", "c"), each = 5), levels = c("a", "b", "c"))
)
tbl %>%
    group_by(y, .drop = FALSE) %>%
    group_rows()
```

 ${\tt group\_by\_drop\_default.DataFrame}$ 

 $Default\ value\ for\ .drop\ argument\ of\ group\_by$ 

#### **Description**

Default value for .drop argument of group\_by

### Usage

```
## S3 method for class 'DataFrame'
group_by_drop_default(.tbl)
```

### **Arguments**

.tbl

A data frame

### Value

TRUE unless .tbl is a grouped data frame that was previously obtained by group\_by(.drop = FALSE)

```
group_by_drop_default(iris)
iris %>%
   group_by(Species) %>%
   group_by_drop_default()

iris %>%
   group_by(Species, .drop = FALSE) %>%
   group_by_drop_default()
```

group\_data 17

grou	n	da	ta

Set and Get Group Data on a DataFrame

#### **Description**

The location of group data is an internal implemnetation detail, so these get and set methods enable interfacing with that data.

### Usage

```
set_group_data(x, g, .drop = group_by_drop_default(x))
get_group_data(x)
```

#### **Arguments**

x A S4Vectors::DataFrame on which to set group data.

g Group data (a data.frame).

.drop Drop groups formed by factor levels that don't appear in the data?

#### Value

For set\_group\_data, the input x with group data set as metadata. For get\_group\_data, the group data that is set on x.

```
group_data.DataFrame Grouping metadata
```

#### **Description**

This collection of functions accesses data about grouped data frames in various ways:

- group\_data() returns a data frame that defines the grouping structure. The columns give the values of the grouping variables. The last column, always called .rows, is a list of integer vectors that gives the location of the rows in each group.
- group\_keys() returns a data frame describing the groups.
- group\_rows() returns a list of integer vectors giving the rows that each group contains.
- group\_indices() returns an integer vector the same length as . data that gives the group that each row belongs to.
- group\_vars() gives names of grouping variables as character vector.
- groups() gives the names of the grouping variables as a list of symbols.
- group\_size() gives the size of each group.
- n\_groups() gives the total number of groups.

See context for equivalent functions that return values for the *current* group.

#### Usage

```
## S3 method for class 'DataFrame'
group_data(.data)
```

### **Arguments**

#### Value

```
a data. frame of group data
```

```
group_vars.DataFrame Grouping metadata
```

### **Description**

This collection of functions accesses data about grouped data frames in various ways:

- group\_data() returns a data frame that defines the grouping structure. The columns give the values of the grouping variables. The last column, always called .rows, is a list of integer vectors that gives the location of the rows in each group.
- group\_keys() returns a data frame describing the groups.
- group\_rows() returns a list of integer vectors giving the rows that each group contains.
- group\_indices() returns an integer vector the same length as .data that gives the group that each row belongs to.
- group\_vars() gives names of grouping variables as character vector.
- groups() gives the names of the grouping variables as a list of symbols.
- group\_size() gives the size of each group.
- n\_groups() gives the total number of groups.

See context for equivalent functions that return values for the *current* group.

### Usage

```
## S3 method for class 'DataFrame'
group_vars(x)
```

### Arguments

```
x a S4Vectors::DataFrame(), likely grouped
```

#### Value

the grouping variables as a character vector

inner\_join.DataFrame 19

```
inner_join.DataFrame Mutating joins
```

### Description

Mutating joins

### Usage

```
## S3 method for class 'DataFrame'
inner_join(
    x,
    y,
    by = NULL,
    copy = FALSE,
    suffix = c(".x", ".y"),
    ...,
    keep = NULL
)
```

### **Arguments**

```
x a DataFrame
y a DataFrame or data.frame
by columns to use for joining the objects. If NULL, the function will look for common columns.
```

#### Value

a DataFrame

```
library(dplyr)
library(S4Vectors)
da <- starwars[, c("name", "mass", "species")][1:10, ]
db <- starwars[, c("name", "homeworld")]

Da <- as(da, "DataFrame")
Db <- as(db, "DataFrame")

Res_inner <- inner_join(Da, Db[1:3, ])</pre>
```

20 mutate.DataFrame

mutate.DataFrame

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

```
## S3 method for class 'DataFrame'
mutate(.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 *Methods*, below, for more details.

...

<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.

#### 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.DataFrame 21

#### **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))
```

22 pull.DataFrame

```
# 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 \leftarrow 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 \leftarrow 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
```

pull.DataFrame

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

```
## S3 method for class 'DataFrame'
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:

• 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 unquote column names and column locations).

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

name

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

```
d <- S4Vectors::DataFrame(mtcars)
pull(d, cyl)</pre>
```

rename, DataFrame-method

Rename Columns of a DataFrame

#### **Description**

Rename Columns of a DataFrame

#### Usage

```
## S4 method for signature 'DataFrame'
rename(x, ...)
```

### **Arguments**

```
x a DataFrame
```

... NSE syntax; new\_name = old\_name to rename selected variables

#### Value

a DataFrame with columns renamed

rename2

Rename Columns of a DataFrame (deprecated)

### **Description**

Rename Columns of a DataFrame (deprecated)

### Usage

```
rename2(.data, ...)
```

### **Arguments**

```
... a DataFrame
... columns to be renamed with syntax new = old
```

#### Value

Deprecated - use rename

select.DataFrame

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

```
## S3 method for class 'DataFrame'
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 *Methods*, below, for more details.
 ... <tidy-select> One or more unquoted expressions separated by commas. Variable names can be used as if they were positions in the data frame, so expressions 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
            3.2
                         1.3
                                     0.2 setosa Sepal.Length
                                                                4.7
#> 4
            3.1
                         1.5
                                     0.2 setosa Sepal.Length
                                                                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
#> 4 Tatooine
                 202
                       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>
#>
                        0.2 setosa Sepal.Length
#> 1
            3.5
                                                 5.1
#> 2
            3.5
                        0.2 setosa Petal.Length
                                                  1.4
#> 3
            3
                        0.2 setosa Sepal.Length
                                                  4.9
#> 4
            3
                        0.2 setosa Petal.Length
                                                  1.4
#> # i 296 more rows
```

#### **Operators::**

The: operator selects a range of consecutive variables:

```
starwars %>% select(name:mass)
```

```
#> # A tibble: 87 x 3
#>
    name
                    height mass
#>
     <chr>
                     <int> <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>
#> 1 blond
                fair
                           blue
                                          19 male masculine Tatooine Human
#> 2 <NA>
                                          112 none masculine Tatooine Droid
                gold
                           yellow
#> 3 <NA>
               white, blue red
                                                none masculine Naboo
                                           33
                                                                          Droid
#> 4 none
               white
                           vellow
                                           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
                         0.2 setosa
             3
#> 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>
                         <dbl> <fct>
#> 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

28 slice.DataFrame

```
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
                          0.2
                                      3.5
#> 2
                          0.2
                                      3
              1.4
#> 3
                                      3.2
              1.3
                          0.2
                                      3.1
#> 4
              1.5
                          0.2
#> # 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()

slice.DataFrame

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

```
## S3 method for class 'DataFrame'
slice(.data, ..., .preserve = FALSE)
```

slice.DataFrame 29

#### **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.

... For slice(): <data-masking> Integer row values.

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.

.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**

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()
```

30 slice.DataFrame

```
# 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)
\ensuremath{\mathtt{\#}} 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, row_number() == n())
filter(mtcars, between(row_number(), 5, n()))
```

summarise.DataFrame 31

summarise.DataFrame

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.

#### Usage

```
## S3 method for class 'DataFrame'
summarise(.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 *Methods*, below, for more details.

<data-masking> Name-value pairs of summary functions. The name will be
the name of the variable in the result.

The value can be:

- A vector of length 1, e.g. min(x), n(), or sum(is.na(y)).
- A data frame, to add multiple columns from a single expression.

[**Deprecated**] Returning values with size 0 or >1 was deprecated as of 1.1.0. Please use reframe() for this instead.

#### 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()
```

32 summarise.DataFrame

#### **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()
```

```
# 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))
```

summarize.DataFrame 33

```
mtcars %>%
   group_by(cyl) %>%
   reframe(qs = quantile(disp, c(0.25, 0.75)), prob = c(0.25, 0.75))
```

summarize.DataFrame

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.
```

#### Usage

```
## S3 method for class 'DataFrame'
summarize(.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 *Methods*, below, for more details.

. . .

<data-masking> Name-value pairs of summary functions. The name will be the name of the variable in the result.

The value can be:

- A vector of length 1, e.g. min(x), n(), or sum(is.na(y)).
- A data frame, to add multiple columns from a single expression.

[**Deprecated**] Returning values with size 0 or >1 was deprecated as of 1.1.0. Please use reframe() for this instead.

#### 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.

34 summarize.DataFrame

#### **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()
```

```
# 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))
```

tally.DataFrame 35

```
# 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))
```

tally.DataFrame

Count the observations in each group

#### **Description**

count() lets you quickly count the unique values of one or more variables: df %>% count(a, b) is roughly equivalent to df %>% group\_by(a, b) %>% summarise(n = n()). count() is paired with tally(), a lower-level helper that is equivalent to df %>% 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

```
## S3 method for class 'DataFrame'
tally(x, wt = NULL, sort = FALSE, name = NULL)
```

#### **Arguments**

x A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g.

from dbplyr or dtplyr).

wt <data-masking> Frequency weights. Can be NULL or a variable:

- 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.

### 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.

36 tbl\_vars.DataFrame

#### **Examples**

```
# 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,
            "male",
  "Max",
                        10,
  "Sandra", "female",
                         1,
  "Susan", "female",
)
# 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)
```

tbl\_vars.DataFrame

List variables provided by a tbl.

### **Description**

tbl\_vars() returns all variables while tbl\_nongroup\_vars() returns only non-grouping variables. The groups attribute of the object returned by tbl\_vars() is a character vector of the grouping columns.

### Usage

```
## S3 method for class 'DataFrame'
tbl_vars(x)
```

ungroup.DataFrame 37

#### **Arguments**

x A tbl object

#### Value

all variables, with a groups attribute when grouped.

#### See Also

group\_vars() for a function that returns grouping variables.

ungroup.DataFrame

Group by one or more variables

### 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

```
## S3 method for class 'DataFrame'
ungroup(x, ...)
```

#### **Arguments**

x A tbl()

... 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.

#### 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.

38 ungroup.DataFrame

#### **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()

```
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
by_vs %>% summarise(n = sum(n))
# To removing grouping, use ungroup
by_vs %>%
  ungroup() %>%
  summarise(n = sum(n))
```

[,DataFrame-method 39

```
# 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(
 x = 1:10,
 y = factor(rep(c("a", "c"), each = 5), levels = c("a", "b", "c"))
tbl %>%
  group_by(y, .drop = FALSE) %>%
  group_rows()
```

[,DataFrame-method

Subset a DataFrame

#### **Description**

Subset a DataFrame

#### Usage

```
## S4 method for signature 'DataFrame'
x[i, j, ..., drop = FALSE]
```

### **Arguments**

```
x a DataFrame to be subset
```

i rows to subset

40 [,DataFrame-method

j columns to subset... other params, passed to regular S4Vectors subsetting

drop drop dimensions?

### Value

a DataFrame subset by rows and/or columns

# Index

* internal group_data, 17	<pre>ends_with(), 24 everything(), 24</pre>
inner_join.DataFrame, 19	f:1+ 4 21 20 20 22 24
+, 20	filter, 4, 21, 28, 29, 32, 34
==, 10	filter(), 29
>, 10	filter.DataFrame, 9
>=, 10	first(), 31, 34
[,DataFrame-method, 39	format.DataFrame, 11
&, <i>10</i>	format.Date, 12
011() 21 24	format.info, 13
all(), 31, 34	format.POSIXct, 12
all_of(), 25	formatC, 13
any(), 31, 34	get_group_data(group_data), 17
any_of(), 25	getOption, 12
arrange, 11, 21, 28, 29, 32, 34	group_by(), 6, 10
arrange(), 7, 14, 15, 38	group_by.DataFrame, 13
arrange.DataFrame, 3	group_by_drop_default(), 14
as.character, 12, 13	group_by_drop_default.DataFrame, 16
AsIs, 12	group_cols(), 24
between(), <i>10</i>	group_data, 17
bindROWS, DataFrame-method, 5	group_data.DataFrame, 17
	group_keys(), 31, 33
$case\_when(), 20$	group_map, 15, 38
cat, <i>12</i>	group_nest, 15, 38
coalesce(), 20	group_split, 15, 38
contains(), 24	group_trim, 15, 38
context, <i>17</i> , <i>18</i>	group_vars(), 37
count.DataFrame, 5	group_vars.DataFrame, 18
cumall(), 20	grouped_df, 14, 28, 31, 33, 37
cumany(), 20	
$cume_dist(), 20$	if_else(), <u>20</u>
cummax(), 20	inner_join.DataFrame, 19
cummean(), 20	IQR(), 31, 34
cummin(), 20	is.na(), <i>10</i>
cumsum(), 20	isS4, <i>12</i>
dense_rank(), 20	lag(), 20
desc, 7	last(), 31, 34
desc(), 4	last_col(), 24
DFplyr (DFplyr-package), 2	lead(), 20
DFplyr-package, 2	$\log(), 20$
distinct.DataFrame, 8	108(), 20
	mad(), <i>31</i> , <i>34</i>
encodeString, <i>13</i>	matches(), 25

42 INDEX

n(), 31, 34 n_distinct(), 31, 34 n_a_if(), 20 nchar, 12 near(), 10 nth(), 31, 34 ntile(), 20 num_range(), 25  paste, 13 percent_rank(), 20 prettyNum, 13 print, 13 print. default, 12 pull. DataFrame, 22  quasiquotation, 23  recode(), 20 reframe, 4, 11, 21, 28, 29, 32, 34 reframe(), 31, 33 rename, 4, 11, 21, 28, 29, 32, 34 rename, DataFrame-method, 23 rename2, 24 row_number(), 20, 29 rowwise, 31, 33  S4Vectors::DataFrame, 17 S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select, 1, 11, 21, 28, 32, 34 select. DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13	mean(), 31, 34 median(), 31, 34 min(), 31, 34 min_rank(), 20 mutate, 4, 11, 28, 29, 32, 34 mutate(), 32, 34	<pre>toString, 12, 13 ungroup(), 10 ungroup.DataFrame, 37 unique.data.frame(), 8 unlist, 12</pre>
percent_rank(), 20 prettyNum, 13 print, 13 print.default, 12 pull.DataFrame, 22  quasiquotation, 23  recode(), 20 reframe, 4, 11, 21, 28, 29, 32, 34 reframe(), 31, 33 rename, 4, 11, 21, 28, 29, 32, 34 rename, DataFrame-method, 23 rename2, 24 row_number(), 20, 29 rowwise, 31, 33  S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select.DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13	<pre>n_distinct(), 31, 34 na_if(), 20 nchar, 12 near(), 10 nth(), 31, 34 ntile(), 20</pre>	
recode(), 20 reframe, 4, 11, 21, 28, 29, 32, 34 reframe(), 31, 33 rename, 4, 11, 21, 28, 29, 32, 34 rename, DataFrame-method, 23 rename2, 24 row_number(), 20, 29 rowwise, 31, 33  S4Vectors::DataFrame, 17 S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select. DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13	percent_rank(), 20 prettyNum, 13 print, 13 print.default, 12	
reframe, 4, 11, 21, 28, 29, 32, 34 reframe(), 31, 33 rename, 4, 11, 21, 28, 29, 32, 34 rename, DataFrame-method, 23 rename2, 24 row_number(), 20, 29 rowwise, 31, 33  S4Vectors::DataFrame, 17 S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select.DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13	quasiquotation, 23	
S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select.DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13	reframe, 4, 11, 21, 28, 29, 32, 34 reframe(), 31, 33 rename, 4, 11, 21, 28, 29, 32, 34 rename, DataFrame-method, 23 rename2, 24 row_number(), 20, 29	
summarise, 4, 11, 21, 28, 29 summarise(), 14, 38 summarise.DataFrame, 31 summarize.DataFrame, 33	S4Vectors::DataFrame(), 18 sd(), 31, 34 select, 4, 11, 21, 29, 32, 34 select.DataFrame, 24 set_group_data (group_data), 17 slice, 4, 11, 21, 28, 32, 34 slice.DataFrame, 28 sprintf, 13 starts_with(), 24, 25 summarise, 4, 11, 21, 28, 29 summarise(), 14, 38 summarise.DataFrame, 31	
<pre>tally.DataFrame, 35 tbl(), 37 tbl_vars.DataFrame, 36</pre>	tbl(), 37	