

# Package ‘tidyr’

January 23, 2020

**Title** Tidy Messy Data

**Version** 1.0.2

**Description** Tools to help to create tidy data, where each column is a variable, each row is an observation, and each cell contains a single value. ‘tidyr’ contains tools for changing the shape (pivoting) and hierarchy (nesting and ‘unnesting’) of a dataset, turning deeply nested lists into rectangular data frames (‘rectangling’), and extracting values out of string columns. It also includes tools for working with missing values (both implicit and explicit).

**License** MIT + file LICENSE

**URL** <https://tidyr.tidyverse.org>,  
<https://github.com/tidyverse/tidyr>

**BugReports** <https://github.com/tidyverse/tidyr/issues>

**Depends** R (>= 3.1)

**Imports** dplyr (>= 0.8.2),  
ellipsis (>= 0.1.0),  
glue,  
magrittr,  
purrr,  
Rcpp,  
rlang,  
stringi,  
tibble (>= 2.1.1),  
tidyselect (>= 0.2.5),  
utils,  
vctrs (>= 0.2.0),  
lifecycle

**Suggests** covr,  
jsonlite,  
knitr,  
repurrrsive (>= 1.0.0),  
rmarkdown,  
readr,  
testthat (>= 2.1.0)

**LinkingTo** Rcpp

**VignetteBuilder** knitr

**Encoding** UTF-8

**LazyData** true

**Roxygen** list(markdown = TRUE)

**RoxygenNote** 7.0.2

## R topics documented:

billboard . . . . .	2
chop . . . . .	3
complete . . . . .	4
construction . . . . .	5
drop_na . . . . .	6
expand . . . . .	7
expand_grid . . . . .	8
extract . . . . .	9
fill . . . . .	10
fish_encounters . . . . .	11
full_seq . . . . .	11
gather . . . . .	12
hoist . . . . .	13
nest . . . . .	17
nest_legacy . . . . .	19
pack . . . . .	21
pivot_longer . . . . .	22
pivot_wider . . . . .	24
relig_income . . . . .	26
replace_na . . . . .	27
separate . . . . .	28
separate_rows . . . . .	29
smiths . . . . .	30
spread . . . . .	31
table1 . . . . .	32
uncount . . . . .	33
unite . . . . .	33
us_rent_income . . . . .	35
who . . . . .	35
world_bank_pop . . . . .	36
<b>Index</b>	<b>37</b>

---

billboard

*Song rankings for billboard top 100 in the year 2000*

---

### Description

Song rankings for billboard top 100 in the year 2000

### Usage

billboard

**Format**

A dataset with variables:

**artist** Artist name

**track** Song name,

**date.enter** Date the song entered the top 100

**wk1 – wk76** Rank of the song in each week after it entered

**Source**

The "Whitburn" project, [https://waxy.org/2008/05/the\\_whitburn\\_project/](https://waxy.org/2008/05/the_whitburn_project/), (downloaded April 2008)

---

 chop

*Chop and unchop*


---

**Description****Maturing**

Chopping and unchopping preserve the width of a data frame, changing its length. `chop()` makes `df` shorter by converting rows within each group into list-columns. `unchop()` makes `df` longer by expanding list-columns so that each element of the list-column gets its own row in the output.

**Usage**

```
chop(data, cols)
```

```
unchop(data, cols, keep_empty = FALSE, ptype = NULL)
```

**Arguments**

<code>data</code>	A data frame.
<code>cols</code>	Column to chop or unchop (automatically quoted). This should be a list-column containing generalised vectors (e.g. any mix of NULLs, atomic vector, S3 vectors, a lists, or data frames).
<code>keep_empty</code>	By default, you get one row of output for each element of the list your unchopping/unnesting. This means that if there's a size-0 element (like NULL or an empty data frame), that entire row will be dropped from the output. If you want to preserve all rows, use <code>keep_empty = TRUE</code> to replace size-0 elements with a single row of missing values.
<code>ptype</code>	Optionally, supply a data frame prototype for the output <code>cols</code> , overriding the default that will be guessed from the combination of individual values.

**Details**

Generally, unchopping is more useful than chopping because it simplifies a complex data structure, and `nest()`ing is usually more appropriate than `chop()`ing since it better preserves the connections between observations.

## Examples

```
# Chop =====
df <- tibble(x = c(1, 1, 1, 2, 2, 3), y = 1:6, z = 6:1)
# Note that we get one row of output for each unique combination of
# non-chopped variables
df %>% chop(c(y, z))
# cf nest
df %>% nest(data = c(y, z))

# Unchop =====
df <- tibble(x = 1:4, y = list(integer(), 1L, 1:2, 1:3))
df %>% unchop(y)
df %>% unchop(y, keep_empty = TRUE)

# Incompatible types -----
# If the list-col contains types that can not be natively
df <- tibble(x = 1:2, y = list("1", 1:3))
try(df %>% unchop(y))
df %>% unchop(y, ptype = tibble(y = integer()))
df %>% unchop(y, ptype = tibble(y = character()))
df %>% unchop(y, ptype = tibble(y = list()))

# Unchopping data frames -----
# Unchopping a list-col of data frames must generate a df-col because
# unchop leaves the column names unchanged
df <- tibble(x = 1:3, y = list(NULL, tibble(x = 1), tibble(y = 1:2)))
df %>% unchop(y)
df %>% unchop(y, keep_empty = TRUE)
```

---

complete

*Complete a data frame with missing combinations of data*

---

## Description

Turns implicit missing values into explicit missing values. This is a wrapper around `expand()`, `dplyr::left_join()` and `replace_na()` that's useful for completing missing combinations of data.

## Usage

```
complete(data, ..., fill = list())
```

## Arguments

`data` A data frame.

`...` Specification of columns to expand. Columns can be atomic vectors or lists. To find all unique combinations of `x`, `y` and `z`, including those not found in the data, supply each variable as a separate argument. To find only the combinations that occur in the data, use `nest: expand(df, nesting(x, y, z))`.

You can combine the two forms. For example, `expand(df, nesting(school_id, student_id), data)` would produce a row for every student for each date.

For factors, the full set of levels (not just those that appear in the data) are used. For continuous variables, you may need to fill in values that don't appear in the data: to do so use expressions like `year = 2010:2020` or `year = \link{full_seq}(year,1)`. Length-zero (empty) elements are automatically dropped.

**fill** A named list that for each variable supplies a single value to use instead of NA for missing combinations.

### Details

If you supply `fill`, these values will also replace existing explicit missing values in the data set.

### Examples

```
library(dplyr, warn.conflicts = FALSE)
df <- tibble(
  group = c(1:2, 1),
  item_id = c(1:2, 2),
  item_name = c("a", "b", "b"),
  value1 = 1:3,
  value2 = 4:6
)
df %>% complete(group, nesting(item_id, item_name))

# You can also choose to fill in missing values
df %>% complete(group, nesting(item_id, item_name), fill = list(value1 = 0))
```

---

construction

*Completed construction in the US in 2018*

---

### Description

Completed construction in the US in 2018

### Usage

```
construction
```

### Format

A dataset with variables:

**Year,Month** Record date

1 unit, 2 to 4 units, 5 units or mote Number of completed units of each size

**Northeast, Midwest, South, West** Number of completed units in each region

### Source

Completions of "New Residential Construction" found in Table 5 at <https://www.census.gov/construction/nrc/xls/newresconst.xls> (downloaded March 2019)

---

drop_na	<i>Drop rows containing missing values</i>
---------	--

---

### Description

Drop rows containing missing values

### Usage

```
drop_na(data, ...)
```

### Arguments

data	A data frame.
...	A selection of columns. If empty, all variables are selected. You can supply bare variable names, select all variables between x and z with <code>x:z</code> , exclude y with <code>-y</code> . For more options, see the <a href="#"><code>dplyr::select()</code></a> documentation. See also the section on selection rules below.

### Rules for selection

Arguments for selecting columns are passed to `tidyselect::vars_select()` and are treated specially. Unlike other verbs, selecting functions make a strict distinction between data expressions and context expressions.

- A data expression is either a bare name like `x` or an expression like `x:y` or `c(x,y)`. In a data expression, you can only refer to columns from the data frame.
- Everything else is a context expression in which you can only refer to objects that you have defined with `<-`.

For instance, `col1:col3` is a data expression that refers to data columns, while `seq(start,end)` is a context expression that refers to objects from the contexts.

If you really need to refer to contextual objects from a data expression, you can unquote them with the tidy eval operator `!!`. This operator evaluates its argument in the context and inlines the result in the surrounding function call. For instance, `c(x, !! x)` selects the `x` column within the data frame and the column referred to by the object `x` defined in the context (which can contain either a column name as string or a column position).

### Examples

```
library(dplyr)
df <- tibble(x = c(1, 2, NA), y = c("a", NA, "b"))
df %>% drop_na()
df %>% drop_na(x)
```

---

expand	<i>Expand data frame to include all combinations of values</i>
--------	--

---

### Description

`expand()` is often useful in conjunction with `left_join()` if you want to convert implicit missing values to explicit missing values. Or you can use it in conjunction with `anti_join()` to figure out which combinations are missing.

### Usage

```
expand(data, ...)
```

```
crossing(...)
```

```
nesting(...)
```

### Arguments

`data` A data frame.

`...` Specification of columns to expand. Columns can be atomic vectors or lists. To find all unique combinations of `x`, `y` and `z`, including those not found in the data, supply each variable as a separate argument. To find only the combinations that occur in the data, use `nest`: `expand(df, nesting(x, y, z))`.

You can combine the two forms. For example, `expand(df, nesting(school_id, student_id), date)` would produce a row for every student for each date.

For factors, the full set of levels (not just those that appear in the data) are used.

For continuous variables, you may need to fill in values that don't appear in the data: to do so use expressions like `year = 2010:2020` or `year = \link{full_seq}(year, 1)`.

Length-zero (empty) elements are automatically dropped.

### Details

`crossing()` is a wrapper around `expand_grid()` that deduplicates and sorts each input. `nesting()` is the complement to `crossing()`: it only keeps combinations of values that appear in the data.

### See Also

`complete()` for a common application of `expand`: completing a data frame with missing combinations. `expand_grid()` is low-level that doesn't deduplicate or sort values.

### Examples

```
library(dplyr)
# All possible combinations of vs & cyl, even those that aren't
# present in the data
expand(mtcars, vs, cyl)
```

```
# Only combinations of vs and cyl that appear in the data
expand(mtcars, nesting(vs, cyl))
```

```

# Implicit missings -----
df <- tibble(
  year = c(2010, 2010, 2010, 2010, 2012, 2012, 2012),
  qtr  = c( 1,  2,  3,  4,  1,  2,  3),
  return = rnorm(7)
)
df %>% expand(year, qtr)
df %>% expand(year = 2010:2012, qtr)
df %>% expand(year = full_seq(year, 1), qtr)
df %>% complete(year = full_seq(year, 1), qtr)

# Nesting -----
# Each person was given one of two treatments, repeated three times
# But some of the replications haven't happened yet, so we have
# incomplete data:
experiment <- tibble(
  name = rep(c("Alex", "Robert", "Sam"), c(3, 2, 1)),
  trt  = rep(c("a", "b", "a"), c(3, 2, 1)),
  rep  = c(1, 2, 3, 1, 2, 1),
  measurement_1 = runif(6),
  measurement_2 = runif(6)
)

# We can figure out the complete set of data with expand()
# Each person only gets one treatment, so we nest name and trt together:
all <- experiment %>% expand(nesting(name, trt), rep)
all

# We can use anti_join to figure out which observations are missing
all %>% anti_join(experiment)

# And use right_join to add in the appropriate missing values to the
# original data
experiment %>% right_join(all)
# Or use the complete() short-hand
experiment %>% complete(nesting(name, trt), rep)

# Generate all combinations with expand():
formulas <- list(
  formula1 = Sepal.Length ~ Sepal.Width,
  formula2 = Sepal.Length ~ Sepal.Width + Petal.Width,
  formula3 = Sepal.Length ~ Sepal.Width + Petal.Width + Petal.Length
)
data <- split(iris, iris$Species)
crossing(formula = formulas, data)

```

---

expand\_grid

*Create a tibble from all combinations of inputs*

---

## Description

Create a tibble from all combinations of inputs

## Usage

```
expand_grid(...)
```



**Arguments**

... Name-value pairs. The name will become the column name in the output.

**Value**

A tibble with one column for each input in ... The output will have one row for each combination of the inputs, i.e. the size be equal to the product of the sizes of the inputs. This implies that if any input has length 0, the output will have zero rows.

**Compared to `expand.grid`**

- Varies the first element fastest.
- Never converts strings to factors.
- Does not add any additional attributes.
- Returns a tibble, not a data frame.
- Can expand any generalised vector, including data frames.

**Examples**

```
expand_grid(x = 1:3, y = 1:2)
expand_grid(l1 = letters, l2 = LETTERS)

# Can also expand data frames
expand_grid(df = data.frame(x = 1:2, y = c(2, 1)), z = 1:3)
# And matrices
expand_grid(x1 = matrix(1:4, nrow = 2), x2 = matrix(5:8, nrow = 2))
```

---

extract	<i>Extract a character column into multiple columns using regular expression groups</i>
---------	---

---

**Description**

Given a regular expression with capturing groups, `extract()` turns each group into a new column. If the groups don't match, or the input is NA, the output will be NA.

**Usage**

```
extract(
  data,
  col,
  into,
  regex = "[[:alnum:]]+",
  remove = TRUE,
  convert = FALSE,
  ...
)
```

**Arguments**

<code>data</code>	A data frame.
<code>col</code>	Column name or position. This is passed to <code>tidyselect::vars_pull()</code> . This argument is passed by expression and supports <a href="#">quasiquote</a> (you can unquote column names or column positions).
<code>into</code>	Names of new variables to create as character vector. Use NA to omit the variable in the output.
<code>regex</code>	a regular expression used to extract the desired values. There should be one group (defined by <code>()</code> ) for each element of <code>into</code> .
<code>remove</code>	If TRUE, remove input column from output data frame.
<code>convert</code>	If TRUE, will run <code>type.convert()</code> with <code>as.is = TRUE</code> on new columns. This is useful if the component columns are integer, numeric or logical. NB: this will cause string "NA"s to be converted to NAs.
<code>...</code>	Additional arguments passed on to methods.

**See Also**

[separate\(\)](#) to split up by a separator.

**Examples**

```
df <- data.frame(x = c(NA, "a-b", "a-d", "b-c", "d-e"))
df %>% extract(x, "A")
df %>% extract(x, c("A", "B"), "([[:alnum:]]+)-([[:alnum:]]+)")

# If no match, NA:
df %>% extract(x, c("A", "B"), "[a-d]+-[a-d]+")
```

---

**fill**

*Fill in missing values with previous or next value*

---

**Description**

Fills missing values in selected columns using the next or previous entry. This is useful in the common output format where values are not repeated, and are only recorded when they change.

**Usage**

```
fill(data, ..., .direction = c("down", "up", "downup", "updown"))
```

**Arguments**

<code>data</code>	A data frame.
<code>...</code>	A selection of columns. If empty, nothing happens. You can supply bare variable names, select all variables between <code>x</code> and <code>z</code> with <code>x:z</code> , exclude <code>y</code> with <code>-y</code> . For more selection options, see the <a href="#">dplyr::select()</a> documentation.
<code>.direction</code>	Direction in which to fill missing values. Currently either "down" (the default), "up", "downup" (i.e. first down and then up) or "updown" (first up and then down).

**Details**

Missing values are replaced in atomic vectors; NULLs are replaced in lists.

**Examples**

```
df <- data.frame(Month = 1:12, Year = c(2000, rep(NA, 11)))
df %>% fill(Year)
```

---

fish_encounters	<i>Fish encounters</i>
-----------------	------------------------

---

**Description**

Information about fish swimming down a river: each station represents an autonomous monitor that records if a tagged fish was seen at that location. Fish travel in one direction (migrating downstream). Information about misses is just as important as hits, but is not directly recorded in this form of the data.

**Usage**

```
fish_encounters
```

**Format**

A dataset with variables:

**fish** Fish identifier

**station** Measurement station

**seen** Was the fish seen? (1 if yes, and true for all rows)

**Source**

Dataset provided by Myfanwy Johnston; more details at <https://fishsciences.github.io/post/visualizing-fish-encounter-histories/>

---

full_seq	<i>Create the full sequence of values in a vector</i>
----------	---

---

**Description**

This is useful if you want to fill in missing values that should have been observed but weren't. For example, `full_seq(c(1, 2, 4, 6), 1)` will return `1:6`.

**Usage**

```
full_seq(x, period, tol = 1e-06)
```

**Arguments**

x	A numeric vector.
period	Gap between each observation. The existing data will be checked to ensure that it is actually of this periodicity.
tol	Numerical tolerance for checking periodicity.

**Examples**

```
full_seq(c(1, 2, 4, 5, 10), 1)
```

---

gather	<i>Gather columns into key-value pairs</i>
--------	--

---

**Description****Retired**

Development on `gather()` is complete, and for new code we recommend switching to `pivot_longer()`, which is easier to use, more featureful, and still under active development. `df %>% gather("key", "value", x, y, z)` is equivalent to `df %>% pivot_longer(c(x, y, z), names_to = "key", values_to = "value")`

See more details in `vignette("pivot")`.

**Usage**

```
gather(
  data,
  key = "key",
  value = "value",
  ...,
  na.rm = FALSE,
  convert = FALSE,
  factor_key = FALSE
)
```

**Arguments**

data	A data frame.
key, value	Names of new key and value columns, as strings or symbols. This argument is passed by expression and supports <a href="#">quasiquote</a> (you can unquote strings and symbols). The name is captured from the expression with <code>rlang::ensym()</code> (note that this kind of interface where symbols do not represent actual objects is now discouraged in the tidyverse; we support it here for backward compatibility).
...	A selection of columns. If empty, all variables are selected. You can supply bare variable names, select all variables between x and z with <code>x:z</code> , exclude y with <code>-y</code> . For more options, see the <code>dplyr::select()</code> documentation. See also the section on selection rules below.
na.rm	If TRUE, will remove rows from output where the value column is NA.
convert	If TRUE will automatically run <code>type.convert()</code> on the key column. This is useful if the column types are actually numeric, integer, or logical.
factor_key	If FALSE, the default, the key values will be stored as a character vector. If TRUE, will be stored as a factor, which preserves the original ordering of the columns.

## Rules for selection

Arguments for selecting columns are passed to `tidyselect::vars_select()` and are treated specially. Unlike other verbs, selecting functions make a strict distinction between data expressions and context expressions.

- A data expression is either a bare name like `x` or an expression like `x:y` or `c(x,y)`. In a data expression, you can only refer to columns from the data frame.
- Everything else is a context expression in which you can only refer to objects that you have defined with `<-`.

For instance, `col1:col3` is a data expression that refers to data columns, while `seq(start,end)` is a context expression that refers to objects from the contexts.

If you really need to refer to contextual objects from a data expression, you can unquote them with the tidy eval operator `!!`. This operator evaluates its argument in the context and inlines the result in the surrounding function call. For instance, `c(x, !! x)` selects the `x` column within the data frame and the column referred to by the object `x` defined in the context (which can contain either a column name as string or a column position).

## Examples

```
library(dplyr)
# From https://stackoverflow.com/questions/1181060
stocks <- tibble(
  time = as.Date('2009-01-01') + 0:9,
  X = rnorm(10, 0, 1),
  Y = rnorm(10, 0, 2),
  Z = rnorm(10, 0, 4)
)

gather(stocks, "stock", "price", -time)
stocks %>% gather("stock", "price", -time)

# get first observation for each Species in iris data -- base R
mini_iris <- iris[c(1, 51, 101), ]
# gather Sepal.Length, Sepal.Width, Petal.Length, Petal.Width
gather(mini_iris, key = "flower_att", value = "measurement",
       Sepal.Length, Sepal.Width, Petal.Length, Petal.Width)
# same result but less verbose
gather(mini_iris, key = "flower_att", value = "measurement", -Species)

# repeat iris example using dplyr and the pipe operator
library(dplyr)
mini_iris <-
  iris %>%
  group_by(Species) %>%
  slice(1)
mini_iris %>% gather(key = "flower_att", value = "measurement", -Species)
```

## Description

### Maturing

`hoist()`, `unnest_longer()`, and `unnest_wider()` provide tools for rectangling, collapsing deeply nested lists into regular columns. `hoist()` allows you to selectively pull components of a list-column out in to their own top-level columns, using the same syntax as `purrr::pluck()`. `unnest_wider()` turns each element of a list-column into a column, and `unnest_longer()` turns each element of a list-column into a row. `unnest_auto()` picks between `unnest_wider()` or `unnest_longer()` based heuristics described below.

Learn more in `vignette("rectangle")`.

## Usage

```
hoist(.data, .col, ..., .remove = TRUE, .simplify = TRUE, .ptype = list())
```

```
unnest_longer(
  data,
  col,
  values_to = NULL,
  indices_to = NULL,
  indices_include = NULL,
  names_repair = "check_unique",
  simplify = TRUE,
  ptype = list()
)
```

```
unnest_wider(
  data,
  col,
  names_sep = NULL,
  simplify = TRUE,
  names_repair = "check_unique",
  ptype = list()
)
```

```
unnest_auto(data, col)
```

## Arguments

<code>.data</code> , <code>data</code>	A data frame.
<code>.col</code> , <code>col</code>	List-column to extract components from.
<code>...</code>	Components of <code>.col</code> to turn into columns in the form <code>col_name = "pluck_specification"</code> . You can pluck by name with a character vector, by position with an integer vector, or with a combination of the two with a list. See <code>purrr::pluck()</code> for details.
<code>.remove</code>	If <code>TRUE</code> , the default, will remove extracted components from <code>.col</code> . This ensures that each value lives only in one place.
<code>.simplify</code>	If <code>TRUE</code> , will attempt to simplify lists of length-1 vectors to an atomic vector
<code>.ptype</code>	Optionally, a named list of prototypes declaring the desired output type of each component.
<code>values_to</code>	Name of column to store vector values. Defaults to <code>col</code> .

<code>indices_to</code>	A string giving the name of column which will contain the inner names or position (if not named) of the values. Defaults to <code>col</code> with <code>_id</code> suffix
<code>indices_include</code>	Add an index column? Defaults to <code>TRUE</code> when <code>col</code> has inner names.
<code>names_repair</code>	Used to check that output data frame has valid names. Must be one of the following options: <ul style="list-style-type: none"> <li>• "minimal": no name repair or checks, beyond basic existence,</li> <li>• "unique": make sure names are unique and not empty,</li> <li>• "check_unique": (the default), no name repair, but check they are unique,</li> <li>• "universal": make the names unique and syntactic</li> <li>• a function: apply custom name repair.</li> <li>• <code>tidyr_legacy</code>: use the name repair from <code>tidyr</code> 0.8.</li> <li>• a formula: a purrr-style anonymous function (see <code>rlang::as_function()</code>)</li> </ul> <p>See <code>vctrs::vec_as_names()</code> for more details on these terms and the strategies used to enforce them.</p>
<code>simplify</code>	If <code>TRUE</code> , will attempt to simplify lists of length-1 vectors to an atomic vector
<code>ptype</code>	Optionally, supply a data frame prototype for the output <code>cols</code> , overriding the default that will be guessed from the combination of individual values.
<code>names_sep</code>	If <code>NULL</code> , the default, the names of new columns will come directly from the inner data frame.  If a string, the names of the new columns will be formed by pasting together the outer column name with the inner names, separated by <code>names_sep</code> .

### Unnest variants

The three `unnest()` functions differ in how they change the shape of the output data frame:

- `unnest_wider()` preserves the rows, but changes the columns.
- `unnest_longer()` preserves the columns, but changes the rows
- `unnest()` can change both rows and columns.

These principles guide their behaviour when they are called with a non-primary data type. For example, if you `unnest_wider()` a list of data frames, the number of rows must be preserved, so each column is turned into a list column of length one. Or if you `unnest_longer()` a list of data frame, the number of columns must be preserved so it creates a packed column. I'm not sure how if these behaviours are useful in practice, but they are theoretically pleasing.

### `unnest_auto()` heuristics

`unnest_auto()` inspects the inner names of the list-col:

- If all elements are unnamed, it uses `unnest_longer()`
- If all elements are named, and there's at least one name in common across all components, it uses `unnest_wider()`
- Otherwise, it falls back to `unnest_longer(indices_include = TRUE)`.

**Examples**

```

df <- tibble(
  character = c("Toothless", "Dory"),
  metadata = list(
    list(
      species = "dragon",
      color = "black",
      films = c(
        "How to Train Your Dragon",
        "How to Train Your Dragon 2",
        "How to Train Your Dragon: The Hidden World"
      )
    ),
    list(
      species = "clownfish",
      color = "blue",
      films = c("Finding Nemo", "Finding Dory")
    )
  )
)
df

# Turn all components of metadata into columns
df %>% unnest_wider(metadata)

# Extract only specified components
df %>% hoist(metadata,
  species = "species",
  first_film = list("films", 1L),
  third_film = list("films", 3L)
)

df %>%
  unnest_wider(metadata) %>%
  unnest_longer(films)
# unnest_longer() is useful when each component of the list should
# form a row
df <- tibble(
  x = 1:3,
  y = list(NULL, 1:3, 4:5)
)
df %>% unnest_longer(y)
# Automatically creates names if widening
df %>% unnest_wider(y)

# And similarly if the vectors are named
df <- tibble(
  x = 1:2,
  y = list(c(a = 1, b = 2), c(a = 10, b = 11, c = 12))
)
df %>% unnest_wider(y)
df %>% unnest_longer(y)

```



nest

*Nest and unnest***Description**

Nesting creates a list-column of data frames; unnesting flattens it back out into regular columns. Nesting is implicitly a summarising operation: you get one row for each group defined by the non-nested columns. This is useful in conjunction with other summaries that work with whole datasets, most notably models.

Learn more in `vignette("nest")`.

**Usage**

```
nest(.data, ..., .key = deprecated())

unnest(
  data,
  cols,
  ...,
  keep_empty = FALSE,
  ptype = NULL,
  names_sep = NULL,
  names_repair = "check_unique",
  .drop = deprecated(),
  .id = deprecated(),
  .sep = deprecated(),
  .preserve = deprecated()
)
```

**Arguments**

<code>.data</code>	A data frame.
<code>...</code>	Name-variable pairs of the form <code>new_col = c(col1, col2, col3)</code> , that describe how you wish to nest existing columns into new columns. The right hand side can be any expression supported by <code>tidyselect</code> . <b>Deprecated:</b> previously you could write <code>df %&gt;% nest(x, y, z)</code> and <code>df %&gt;% unnest(x, y, z)</code> . Convert to <code>df %&gt;% nest(data = c(x, y, z))</code> . and <code>df %&gt;% unnest(c(x, y, z))</code> . If you previously created new variable in <code>unnest()</code> you'll now need to do it explicitly with <code>mutate()</code> . Convert <code>df %&gt;% unnest(y = fun(x, y, z))</code> to <code>df %&gt;% mutate(y = fun(x, y, z)) %&gt;% unnest(y)</code> .
<code>.key</code>	<b>Deprecated:</b> No longer needed because of the new <code>new_col = c(col1, col2, col3)</code> syntax.
<code>data</code>	A data frame.
<code>cols</code>	Names of columns to unnest. If you <code>unnest()</code> multiple columns, parallel entries must compatible sizes, i.e. they're either equal or length 1 (following the standard tidyverse recycling rules).
<code>keep_empty</code>	By default, you get one row of output for each element of the list your unchopping/unnesting. This means that if there's a size-0 element (like <code>NULL</code> or an empty data frame), that entire row will be dropped from the output. If you want

	to preserve all rows, use <code>keep_empty = TRUE</code> to replace size-0 elements with a single row of missing values.
<code>ptype</code>	Optionally, supply a data frame prototype for the output <code>cols</code> , overriding the default that will be guessed from the combination of individual values.
<code>names_sep</code>	If <code>NULL</code> , the default, the names of new columns will come directly from the inner data frame. If a string, the names of the new columns will be formed by pasting together the outer column name with the inner names, separated by <code>names_sep</code> .
<code>names_repair</code>	Used to check that output data frame has valid names. Must be one of the following options: <ul style="list-style-type: none"> <li>• "minimal": no name repair or checks, beyond basic existence,</li> <li>• "unique": make sure names are unique and not empty,</li> <li>• "check_unique": (the default), no name repair, but check they are unique,</li> <li>• "universal": make the names unique and syntactic</li> <li>• a function: apply custom name repair.</li> <li>• <code>tidyr_legacy</code>: use the name repair from tidyr 0.8.</li> <li>• a formula: a purrr-style anonymous function (see <a href="#">rlang::as_function()</a>)</li> </ul> See <a href="#">vctrs::vec_as_names()</a> for more details on these terms and the strategies used to enforce them.
<code>.drop, .preserve</code>	<b>Deprecated:</b> all list-columns are now preserved; If there are any that you don't want in the output use <code>select()</code> to remove them prior to unnesting.
<code>.id</code>	<b>Deprecated:</b> convert <code>df %&gt;% unnest(x, .id = "id")</code> to <code>df %&gt;% mutate(id = names(x)) %&gt;% unnest(x)</code>
<code>.sep</code>	<b>Deprecated:</b> use <code>names_sep</code> instead.

### New syntax

tidyr 1.0.0 introduced a new syntax for `nest()` and `unnest()` that's designed to be more similar to other functions. Converting to the new syntax should be straightforward (guided by the message you'll receive) but if you just need to run an old analysis, you can easily revert to the previous behaviour using [nest\\_legacy\(\)](#) and [unnest\\_legacy\(\)](#) as follows:

```
library(tidyr)
nest <- nest_legacy
unnest <- unnest_legacy
```

### Grouped data frames

`df %>% nest(x, y)` specifies the columns to be nested; i.e. the columns that will appear in the inner data frame. Alternatively, you can `nest()` a grouped data frame created by [dplyr::group\\_by\(\)](#). The grouping variables remain in the outer data frame and the others are nested. The result preserves the grouping of the input.

Variables supplied to `nest()` will override grouping variables so that `df %>% group_by(x, y) %>% nest(z)` will be equivalent to `df %>% nest(z)`.

### Examples

```
df <- tibble(x = c(1, 1, 1, 2, 2, 3), y = 1:6, z = 6:1)
# Note that we get one row of output for each unique combination of
# non-nested variables
```

```

df %>% nest(data = c(y, z))
# chop does something similar, but retains individual columns
df %>% chop(c(y, z))

# use tidysselect syntax and helpers, just like in dplyr::select()
df %>% nest(data = one_of("y", "z"))

iris %>% nest(data = -Species)
nest_vars <- names(iris)[1:4]
iris %>% nest(data = one_of(nest_vars))
iris %>%
  nest(petal = starts_with("Petal"), sepal = starts_with("Sepal"))
iris %>%
  nest(width = contains("Width"), length = contains("Length"))

# Nesting a grouped data frame nests all variables apart from the group vars
library(dplyr)
fish_encounters %>%
  group_by(fish) %>%
  nest()

# Nesting is often useful for creating per group models
mtcars %>%
  group_by(cyl) %>%
  nest() %>%
  mutate(models = lapply(data, function(df) lm(mpg ~ wt, data = df)))

# unnest() is primarily designed to work with lists of data frames
df <- tibble(
  x = 1:3,
  y = list(
    NULL,
    tibble(a = 1, b = 2),
    tibble(a = 1:3, b = 3:1)
  )
)
df %>% unnest(y)
df %>% unnest(y, keep_empty = TRUE)

# If you have lists of lists, or lists of atomic vectors, instead
# see hoist(), unnest_wider(), and unnest_longer()

#' # You can unnest multiple columns simultaneously
df <- tibble(
  a = list(c("a", "b"), "c"),
  b = list(1:2, 3),
  c = c(11, 22)
)
df %>% unnest(c(a, b))

# Compare with unnesting one column at a time, which generates
# the Cartesian product
df %>% unnest(a) %>% unnest(b)

```

## Description

### Retired

tidyr 1.0.0 introduced a new syntax for `nest()` and `unnest()`. The majority of existing usage should be automatically translated to the new syntax with a warning. However, if you need to quickly roll back to the previous behaviour, these functions provide the previous interface. To make old code work as is, add the following code to the top of your script:

```
library(tidyr)
nest <- nest_legacy
unnest <- unnest_legacy
```

## Usage

```
nest_legacy(data, ..., .key = "data")

unnest_legacy(data, ..., .drop = NA, .id = NULL, .sep = NULL, .preserve = NULL)
```

## Arguments

<code>data</code>	A data frame.
<code>...</code>	Specification of columns to unnest. Use bare variable names or functions of variables. If omitted, defaults to all list-cols.
<code>.key</code>	The name of the new column, as a string or symbol. This argument is passed by expression and supports <a href="#">quasiquote</a> (you can unquote strings and symbols). The name is captured from the expression with <code>rlang::ensym()</code> (note that this kind of interface where symbols do not represent actual objects is now discouraged in the tidyverse; we support it here for backward compatibility).
<code>.drop</code>	Should additional list columns be dropped? By default, <code>unnest()</code> will drop them if unnesting the specified columns requires the rows to be duplicated.
<code>.id</code>	Data frame identifier - if supplied, will create a new column with name <code>.id</code> , giving a unique identifier. This is most useful if the list column is named.
<code>.sep</code>	If non-NULL, the names of unnested data frame columns will combine the name of the original list-col with the names from the nested data frame, separated by <code>.sep</code> .
<code>.preserve</code>	Optionally, list-columns to preserve in the output. These will be duplicated in the same way as atomic vectors. This has <code>dplyr::select()</code> semantics so you can preserve multiple variables with <code>.preserve = c(x,y)</code> or <code>.preserve = starts_with("list")</code> .

## Examples

```
# Nest and unnest are inverses
df <- data.frame(x = c(1, 1, 2), y = 3:1)
df %>% nest_legacy(y)
df %>% nest_legacy(y) %>% unnest_legacy()

# nesting -----
as_tibble(iris) %>% nest_legacy(-Species)
as_tibble(chickwts) %>% nest_legacy(weight)

# unnesting -----
```

```
df <- tibble(
  x = 1:2,
  y = list(
    tibble(z = 1),
    tibble(z = 3:4)
  )
)
df %>% unnest_legacy(y)

# You can also unnest multiple columns simultaneously
df <- tibble(
  a = list(c("a", "b"), "c"),
  b = list(1:2, 3),
  c = c(11, 22)
)
df %>% unnest_legacy(a, b)
# If you omit the column names, it'll unnest all list-cols
df %>% unnest_legacy()
```

---

pack

*Pack and unpack*


---

## Description

### Maturing

Packing and unpacking preserve the length of a data frame, changing its width. `pack()` makes df narrow by collapsing a set of columns into a single df-column. `unpack()` makes data wider by expanding df-columns back out into individual columns.

## Usage

```
pack(data, ...)
```

```
unpack(data, cols, names_sep = NULL, names_repair = "check_unique")
```

## Arguments

<code>data</code>	A data frame.
<code>...</code>	Name-variable pairs of the form <code>new_col = c(col1, col2, col3)</code> , that describe how you wish to pack existing columns into new columns. The right hand side can be any expression supported by <code>tidyselect</code> .
<code>cols</code>	Name of column that you wish to unpack.
<code>names_sep</code>	If <code>NULL</code> , the default, the names of new columns will come directly from the inner data frame. If a string, the names of the new columns will be formed by pasting together the outer column name with the inner names, separated by <code>names_sep</code> .
<code>names_repair</code>	Used to check that output data frame has valid names. Must be one of the following options: <ul style="list-style-type: none"> <li>"minimal": no name repair or checks, beyond basic existence,</li> <li>"unique": make sure names are unique and not empty,</li> </ul>

- "check\_unique": (the default), no name repair, but check they are unique,
- "universal": make the names unique and syntactic
- a function: apply custom name repair.
- `tidyr_legacy`: use the name repair from tidyr 0.8.
- a formula: a purrr-style anonymous function (see `rlang::as_function()`)

See `vctrs::vec_as_names()` for more details on these terms and the strategies used to enforce them.

## Details

Generally, unpacking is more useful than packing because it simplifies a complex data structure. Currently, few functions work with df-cols, and they are mostly a curiosity, but seem worth exploring further because they mimic the nested column headers that are so popular in Excel.

## Examples

```
# Packing =====
# It's not currently clear why you would ever want to pack columns
# since few functions work with this sort of data.
df <- tibble(x1 = 1:3, x2 = 4:6, x3 = 7:9, y = 1:3)
df
df %>% pack(x = starts_with("x"))
df %>% pack(x = c(x1, x2, x3), y = y)

# Unpacking =====
df <- tibble(
  x = 1:3,
  y = tibble(a = 1:3, b = 3:1),
  z = tibble(X = c("a", "b", "c"), Y = runif(3), Z = c(TRUE, FALSE, NA))
)
df
df %>% unpack(y)
df %>% unpack(c(y, z))
df %>% unpack(c(y, z), names_sep = "_")
```

---

pivot\_longer

*Pivot data from wide to long*

---

## Description

### Maturing

`pivot_longer()` "lengthens" data, increasing the number of rows and decreasing the number of columns. The inverse transformation is `pivot_wider()`

Learn more in `vignette("pivot")`.

## Usage

```
pivot_longer(
  data,
  cols,
  names_to = "name",
```

```

names_prefix = NULL,
names_sep = NULL,
names_pattern = NULL,
names_ptypes = list(),
names_repair = "check_unique",
values_to = "value",
values_drop_na = FALSE,
values_ptypes = list()
)

```

## Arguments

<code>data</code>	A data frame to pivot.
<code>cols</code>	Columns to pivot into longer format. This takes a tidysselect specification.
<code>names_to</code>	A string specifying the name of the column to create from the data stored in the column names of data. Can be a character vector, creating multiple columns, if <code>names_sep</code> or <code>names_pattern</code> is provided.
<code>names_prefix</code>	A regular expression used to remove matching text from the start of each variable name.
<code>names_sep, names_pattern</code>	If <code>names_to</code> contains multiple values, these arguments control how the column name is broken up. <code>names_sep</code> takes the same specification as <code>separate()</code> , and can either be a numeric vector (specifying positions to break on), or a single string (specifying a regular expression to split on). <code>names_pattern</code> takes the same specification as <code>extract()</code> , a regular expression containing matching groups ( <code>()</code> ). If these arguments does not give you enough control, use <code>pivot_longer_spec()</code> to create a spec object and process manually as needed.
<code>names_ptypes, values_ptypes</code>	A list of of column name-prototype pairs. A prototype (or ptype for short) is a zero-length vector (like <code>integer()</code> or <code>numeric()</code> ) that defines the type, class, and attributes of a vector. If not specified, the type of the columns generated from <code>names_to</code> will be character, and the type of the variables generated from <code>values_to</code> will be the common type of the input columns used to generate them.
<code>names_repair</code>	What happen if the output has invalid column names? The default, "check_unique" is to error if the columns are duplicated. Use "minimal" to allow duplicates in the output, or "unique" to de-duplicated by adding numeric suffixes. See <code>vctrs::vec_as_names()</code> for more options.
<code>values_to</code>	A string specifying the name of the column to create from the data stored in cell values. If <code>names_to</code> is a character containing the special <code>.value</code> sentinel, this value will be ignored, and the name of the value column will be derived from part of the existing column names.
<code>values_drop_na</code>	If TRUE, will drop rows that contain only NAs in the <code>value_to</code> column. This effectively converts explicit missing values to implicit missing values, and should generally be used only when missing values in data were created by its structure.

## Details

`pivot_longer()` is an updated approach to `gather()`, designed to be both simpler to use and to handle more use cases. We recommend you use `pivot_longer()` for new code; `gather()` isn't going away but is no longer under active development.

## Examples

```
# See vignette("pivot") for examples and explanation

# Simplest case where column names are character data
relig_income
relig_income %>%
  pivot_longer(-religion, names_to = "income", values_to = "count")

# Slightly more complex case where columns have common prefix,
# and missing missings are structural so should be dropped.
billboard
billboard %>%
  pivot_longer(
    cols = starts_with("wk"),
    names_to = "week",
    names_prefix = "wk",
    values_to = "rank",
    values_drop_na = TRUE
  )

# Multiple variables stored in column names
who %>% pivot_longer(
  cols = new_sp_m014:newrel_f65,
  names_to = c("diagnosis", "gender", "age"),
  names_pattern = "new_?(.*)_(.)(.*)",
  values_to = "count"
)

# Multiple observations per row
anscombe
anscombe %>%
  pivot_longer(everything(),
    names_to = c(".value", "set"),
    names_pattern = "(.)(.)(.)"
  )
```

---

pivot\_wider

*Pivot data from long to wide*

---

## Description

### Maturing

`pivot_wider()` "widens" data, increasing the number of columns and decreasing the number of rows. The inverse transformation is `pivot_longer()`.

Learn more in `vignette("pivot")`.



**Usage**

```

pivot_wider(
  data,
  id_cols = NULL,
  names_from = name,
  names_prefix = "",
  names_sep = "_",
  names_repair = "check_unique",
  values_from = value,
  values_fill = NULL,
  values_fn = NULL
)

```

**Arguments**

<code>data</code>	A data frame to pivot.
<code>id_cols</code>	A set of columns that uniquely identifies each observation. Defaults to all columns in <code>data</code> except for the columns specified in <code>names_from</code> and <code>values_from</code> . Typically used when you have additional variables that is directly related.
<code>names_from</code> , <code>values_from</code>	A pair of arguments describing which column (or columns) to get the name of the output column ( <code>names_from</code> ), and which column (or columns) to get the cell values from ( <code>values_from</code> ). If <code>values_from</code> contains multiple values, the value will be added to the front of the output column.
<code>names_prefix</code>	String added to the start of every variable name. This is particularly useful if <code>names_from</code> is a numeric vector and you want to create syntactic variable names.
<code>names_sep</code>	If <code>names_from</code> or <code>values_from</code> contains multiple variables, this will be used to join their values together into a single string to use as a column name.
<code>names_repair</code>	What happen if the output has invalid column names? The default, "check_unique" is to error if the columns are duplicated. Use "minimal" to allow duplicates in the output, or "unique" to de-duplicated by adding numeric suffixes. See <a href="#">vctrs::vec_as_names()</a> for more options.
<code>values_fill</code>	Optionally, a named list specifying what each value should be filled in with when missing.
<code>values_fn</code>	Optionally, a named list providing a function that will be applied to the value in each cell in the output. You will typically use this when the combination of <code>id_cols</code> and <code>value</code> column does not uniquely identify an observation.

**Details**

`pivot_wider()` is an updated approach to [spread\(\)](#), designed to be both simpler to use and to handle more use cases. We recommend you use `pivot_wider()` for new code; `spread()` isn't going away but is no longer under active development.

**See Also**

[pivot\\_wider\\_spec\(\)](#) to pivot "by hand" with a data frame that defines a pivoting specification.

**Examples**

```

# See vignette("pivot") for examples and explanation

fish_encounters
fish_encounters %>%
  pivot_wider(names_from = station, values_from = seen)
# Fill in missing values
fish_encounters %>%
  pivot_wider(
    names_from = station,
    values_from = seen,
    values_fill = list(seen = 0)
  )

# Generate column names from multiple variables
us_rent_income %>%
  pivot_wider(names_from = variable, values_from = c(estimate, moe))

# Can perform aggregation with values_fn
warpbreaks <- as_tibble(warpbreaks[c("wool", "tension", "breaks")])
warpbreaks
warpbreaks %>%
  pivot_wider(
    names_from = wool,
    values_from = breaks,
    values_fn = list(breaks = mean)
  )

```

---

relig\_income

*Pew religion and income survey*


---

**Description**

Pew religion and income survey

**Usage**

```
relig_income
```

**Format**

A dataset with variables:

**religion** Name of religion

<\$10k-Don't know/refused Number of respondees with income range in column name

**Source**

Downloaded from <https://www.pewforum.org/religious-landscape-study/> (downloaded November 2009)

---

replace_na	<i>Replace missing values</i>
------------	-------------------------------

---

### Description

Replace missing values

### Usage

```
replace_na(data, replace, ...)
```

### Arguments

data	A data frame or vector.
replace	If data is a data frame, a named list giving the value to replace NA with for each column. If data is a vector, a single value used for replacement.
...	Additional arguments for methods. Currently unused.

### Value

If data is a data frame, returns a data frame. If data is a vector, returns a vector of class determined by the union of data and replace.

### See Also

[na\\_if](#) to replace specified values with a NA. [coalesce](#) to replace missing values with a specified value. [recode](#) to more generally replace values.

### Examples

```
library(dplyr)
df <- tibble(x = c(1, 2, NA), y = c("a", NA, "b"), z = list(1:5, NULL, 10:20))
df %>% replace_na(list(x = 0, y = "unknown"))
df %>% mutate(x = replace_na(x, 0))

# NULL are the list-col equivalent of NAs
df %>% replace_na(list(z = list(5)))

df$x %>% replace_na(0)
df$y %>% replace_na("unknown")
```

---

separate	<i>Separate a character column into multiple columns using a regular expression separator</i>
----------	---

---

### Description

Given either regular expression or a vector of character positions, `separate()` turns a single character column into multiple columns.

### Usage

```
separate(
  data,
  col,
  into,
  sep = "[^[:alnum:]]+",
  remove = TRUE,
  convert = FALSE,
  extra = "warn",
  fill = "warn",
  ...
)
```

### Arguments

data	A data frame.
col	Column name or position. This is passed to <code>tidyselect::vars_pull()</code> . This argument is passed by expression and supports <a href="#">quasiquotation</a> (you can unquote column names or column positions).
into	Names of new variables to create as character vector. Use NA to omit the variable in the output.
sep	Separator between columns. If character, is interpreted as a regular expression. The default value is a regular expression that matches any sequence of non-alphanumeric values. If numeric, interpreted as positions to split at. Positive values start at 1 at the far-left of the string; negative value start at -1 at the far-right of the string. The length of sep should be one less than into.
remove	If TRUE, remove input column from output data frame.
convert	If TRUE, will run <code>type.convert()</code> with <code>as.is = TRUE</code> on new columns. This is useful if the component columns are integer, numeric or logical. NB: this will cause string "NA"s to be converted to NAs.
extra	If sep is a character vector, this controls what happens when there are too many pieces. There are three valid options: <ul style="list-style-type: none"> <li>• "warn" (the default): emit a warning and drop extra values.</li> <li>• "drop": drop any extra values without a warning.</li> <li>• "merge": only splits at most <code>length(into)</code> times</li> </ul>
fill	If sep is a character vector, this controls what happens when there are not enough pieces. There are three valid options:

- "warn" (the default): emit a warning and fill from the right
  - "right": fill with missing values on the right
  - "left": fill with missing values on the left
- ... Additional arguments passed on to methods.

### See Also

`unite()`, the complement, `extract()` which uses regular expression capturing groups.

### Examples

```
library(dplyr)
df <- data.frame(x = c(NA, "a.b", "a.d", "b.c"))
df %>% separate(x, c("A", "B"))

# If you just want the second variable:
df %>% separate(x, c(NA, "B"))

# If every row doesn't split into the same number of pieces, use
# the extra and fill arguments to control what happens
df <- data.frame(x = c("a", "a b", "a b c", NA))
df %>% separate(x, c("a", "b"))
# The same behaviour drops the c but no warnings
df %>% separate(x, c("a", "b"), extra = "drop", fill = "right")
# Another option:
df %>% separate(x, c("a", "b"), extra = "merge", fill = "left")
# Or you can keep all three
df %>% separate(x, c("a", "b", "c"))

# If only want to split specified number of times use extra = "merge"
df <- data.frame(x = c("x: 123", "y: error: 7"))
df %>% separate(x, c("key", "value"), ": ", extra = "merge")

# Use regular expressions to separate on multiple characters:
df <- data.frame(x = c(NA, "a?b", "a.d", "b:c"))
df %>% separate(x, c("A", "B"), sep = "[\\.\\?\\:]", fill = "right")

# convert = TRUE detects column classes
df <- data.frame(x = c("a:1", "a:2", "c:4", "d", NA))
df %>% separate(x, c("key", "value"), ":") %>% str
df %>% separate(x, c("key", "value"), ":", convert = TRUE) %>% str

# Argument col can take quasiquotation to work with strings
var <- "x"
df %>% separate(!var, c("key", "value"), ":")
```

---

separate\_rows

*Separate a collapsed column into multiple rows*

---

### Description

If a variable contains observations with multiple delimited values, this separates the values and places each one in its own row.

**Usage**

```
separate_rows(data, ..., sep = "[^[:alnum:]]+", convert = FALSE)
```

**Arguments**

data	A data frame.
...	A selection of columns. If empty, nothing happens. You can supply bare variable names, select all variables between x and z with x:z, exclude y with -y. For more selection options, see the <code>dplyr::select()</code> documentation.
sep	Separator delimiting collapsed values.
convert	If TRUE will automatically run <code>type.convert()</code> on the key column. This is useful if the column types are actually numeric, integer, or logical.

**Rules for selection**

Arguments for selecting columns are passed to `tidyselect::vars_select()` and are treated specially. Unlike other verbs, selecting functions make a strict distinction between data expressions and context expressions.

- A data expression is either a bare name like `x` or an expression like `x:y` or `c(x,y)`. In a data expression, you can only refer to columns from the data frame.
- Everything else is a context expression in which you can only refer to objects that you have defined with `<-`.

For instance, `col1:col3` is a data expression that refers to data columns, while `seq(start, end)` is a context expression that refers to objects from the contexts.

If you really need to refer to contextual objects from a data expression, you can unquote them with the tidy eval operator `!!`. This operator evaluates its argument in the context and inlines the result in the surrounding function call. For instance, `c(x, !! x)` selects the `x` column within the data frame and the column referred to by the object `x` defined in the context (which can contain either a column name as string or a column position).

**Examples**

```
df <- data.frame(
  x = 1:3,
  y = c("a", "d,e,f", "g,h"),
  z = c("1", "2,3,4", "5,6"),
  stringsAsFactors = FALSE
)
separate_rows(df, y, z, convert = TRUE)
```

---

 smiths

*Some data about the Smith family*


---

**Description**

A small demo dataset describing John and Mary Smith.

**Usage**

```
smiths
```

**Format**

A data frame with 2 rows and 5 columns.

---

spread	<i>Spread a key-value pair across multiple columns</i>
--------	--

---

**Description****Retired**

Development on `spread()` is complete, and for new code we recommend switching to `pivot_wider()`, which is easier to use, more featureful, and still under active development. `df %>% spread(key, value)` is equivalent to `df %>% pivot_wider(names_from = key, values_from = value)`

See more details in `vignette("pivot")`.

**Usage**

```
spread(data, key, value, fill = NA, convert = FALSE, drop = TRUE, sep = NULL)
```

**Arguments**

<code>data</code>	A data frame.
<code>key, value</code>	Column names or positions. This is passed to <code>tidyselect::vars_pull()</code> . These arguments are passed by expression and support <a href="#">quasiquotation</a> (you can unquote column names or column positions).
<code>fill</code>	If set, missing values will be replaced with this value. Note that there are two types of missingness in the input: explicit missing values (i.e. <code>NA</code> ), and implicit missings, rows that simply aren't present. Both types of missing value will be replaced by <code>fill</code> .
<code>convert</code>	If <code>TRUE</code> , <code>type.convert()</code> with <code>asis = TRUE</code> will be run on each of the new columns. This is useful if the value column was a mix of variables that was coerced to a string. If the class of the value column was factor or date, note that will not be true of the new columns that are produced, which are coerced to character before type conversion.
<code>drop</code>	If <code>FALSE</code> , will keep factor levels that don't appear in the data, filling in missing combinations with <code>fill</code> .
<code>sep</code>	If <code>NULL</code> , the column names will be taken from the values of key variable. If non- <code>NULL</code> , the column names will be given by " <code>&lt;key_name&gt;&lt;sep&gt;&lt;key_value&gt;</code> ".

**Examples**

```
library(dplyr)
stocks <- data.frame(
  time = as.Date('2009-01-01') + 0:9,
  X = rnorm(10, 0, 1),
  Y = rnorm(10, 0, 2),
  Z = rnorm(10, 0, 4)
)
stocksm <- stocks %>% gather(stock, price, -time)
stocksm %>% spread(stock, price)
```

```

stocksm %>% spread(time, price)

# Spread and gather are complements
df <- data.frame(x = c("a", "b"), y = c(3, 4), z = c(5, 6))
df %>% spread(x, y) %>% gather("x", "y", a:b, na.rm = TRUE)

# Use 'convert = TRUE' to produce variables of mixed type
df <- data.frame(row = rep(c(1, 51), each = 3),
                 var = c("Sepal.Length", "Species", "Species_num"),
                 value = c(5.1, "setosa", 1, 7.0, "versicolor", 2))
df %>% spread(var, value) %>% str
df %>% spread(var, value, convert = TRUE) %>% str

```

table1

*Example tabular representations***Description**

Data sets that demonstrate multiple ways to layout the same tabular data.

**Usage**

table1

table2

table3

table4a

table4b

table5

**Format**

An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 6 rows and 4 columns.

**Details**

table1, table2, table3, table4a, table4b, and table5 all display the number of TB cases documented by the World Health Organization in Afghanistan, Brazil, and China between 1999 and 2000. The data contains values associated with four variables (country, year, cases, and population), but each table organizes the values in a different layout.

The data is a subset of the data contained in the World Health Organization Global Tuberculosis Report

**Source**

<https://www.who.int/tb/country/data/download/en/>



---

uncount	<i>"Uncount" a data frame</i>
---------	-------------------------------

---

### Description

Performs the opposite operation to `dplyr::count()`, duplicating rows according to a weighting variable (or expression).

### Usage

```
uncount(data, weights, .remove = TRUE, .id = NULL)
```

### Arguments

<code>data</code>	A data frame, tibble, or grouped tibble.
<code>weights</code>	A vector of weights. Evaluated in the context of data; supports quasiquotation.
<code>.remove</code>	If TRUE, and <code>weights</code> is a single
<code>.id</code>	Supply a string to create a new variable which gives a unique identifier for each created row.

### Examples

```
df <- tibble(x = c("a", "b"), n = c(1, 2))
uncount(df, n)
uncount(df, n, .id = "id")

# You can also use constants
uncount(df, 2)

# Or expressions
uncount(df, 2 / n)
```

---

unite	<i>Unite multiple columns into one by pasting strings together</i>
-------	--

---

### Description

Convenience function to paste together multiple columns into one.

### Usage

```
unite(data, col, ..., sep = "_", remove = TRUE, na.rm = FALSE)
```

**Arguments**

<code>data</code>	A data frame.
<code>col</code>	The name of the new column, as a string or symbol. This argument is passed by expression and supports <a href="#">quasiquote</a> (you can unquote strings and symbols). The name is captured from the expression with <code>rlang::ensym()</code> (note that this kind of interface where symbols do not represent actual objects is now discouraged in the tidyverse; we support it here for backward compatibility).
<code>...</code>	A selection of columns. If empty, all variables are selected. You can supply bare variable names, select all variables between <code>x</code> and <code>z</code> with <code>x:z</code> , exclude <code>y</code> with <code>-y</code> . For more options, see the <code>dplyr::select()</code> documentation. See also the section on selection rules below.
<code>sep</code>	Separator to use between values.
<code>remove</code>	If TRUE, remove input columns from output data frame.
<code>na.rm</code>	If TRUE, missing values will be removed prior to uniting each value.

**Rules for selection**

Arguments for selecting columns are passed to `tidyselect::vars_select()` and are treated specially. Unlike other verbs, selecting functions make a strict distinction between data expressions and context expressions.

- A data expression is either a bare name like `x` or an expression like `x:y` or `c(x,y)`. In a data expression, you can only refer to columns from the data frame.
- Everything else is a context expression in which you can only refer to objects that you have defined with `<-`.

For instance, `col1:col3` is a data expression that refers to data columns, while `seq(start,end)` is a context expression that refers to objects from the contexts.

If you really need to refer to contextual objects from a data expression, you can unquote them with the tidy eval operator `!!`. This operator evaluates its argument in the context and inlines the result in the surrounding function call. For instance, `c(x, !! x)` selects the `x` column within the data frame and the column referred to by the object `x` defined in the context (which can contain either a column name as string or a column position).

**See Also**

[separate\(\)](#), the complement.

**Examples**

```
df <- expand_grid(x = c("a", NA), y = c("b", NA))
df

df %>% unite("z", x:y, remove = FALSE)
# To remove missing values:
df %>% unite("z", x:y, na.rm = TRUE, remove = FALSE)

# Separate is almost the complement of unite
df %>%
  unite("xy", x:y) %>%
  separate(xy, c("x", "y"))
# (but note `x` and `y` contain now "NA" not NA)
```

---

us_rent_income	<i>US rent and income data</i>
----------------	--------------------------------

---

**Description**

Captured from the 2017 American Community Survey using the tidycensus package.

**Usage**

us\_rent\_income

**Format**

A dataset with variables:

**GEOID** FIP state identifier

**NAME** Name of state

**variable** Variable name: income = median yearly income, rent = median monthly rent

**estimate** Estimated value

**moe** 90% margin of error

---

who	<i>World Health Organization TB data</i>
-----	--

---

**Description**

A subset of data from the World Health Organization Global Tuberculosis Report, and accompanying global populations.

**Usage**

who

population

**Format**

A dataset with the variables

**country** Country name

**iso2, iso3** 2 & 3 letter ISO country codes

**year** Year

**new\_sp\_m014 - new\_rel\_f65** Counts of new TB cases recorded by group. Column names encode three variables that describe the group (see details).

**Details**

The data uses the original codes given by the World Health Organization. The column names for columns five through 60 are made by combining new\_ to a code for method of diagnosis (rel = relapse, sn = negative pulmonary smear, sp = positive pulmonary smear, ep = extrapulmonary) to a code for gender (f = female, m = male) to a code for age group (014 = 0-14 yrs of age, 1524 = 15-24 years of age, 2534 = 25 to 34 years of age, 3544 = 35 to 44 years of age, 4554 = 45 to 54 years of age, 5564 = 55 to 64 years of age, 65 = 65 years of age or older).

**Source**

<https://www.who.int/tb/country/data/download/en/>

---

world\_bank\_pop

*Population data from the world bank*

---

**Description**

Data about population from the World Bank.

**Usage**

world\_bank\_pop

**Format**

A dataset with variables:

**country** Three letter country code

**indicator** Indicator name: SP.POP.GROW = population growth, SP.POP.TOTL = total population, SP.URB.GROW = urban population growth, SP.URB.TOTL = total urban population

**2000-2018** Value for each year

**Source**

Dataset from the World Bank data bank: <https://data.worldbank.org>

# Index

## \*Topic **datasets**

- billboard, [2](#)
  - construction, [5](#)
  - fish\_encounters, [11](#)
  - relig\_income, [26](#)
  - smiths, [30](#)
  - table1, [32](#)
  - us\_rent\_income, [35](#)
  - who, [35](#)
  - world\_bank\_pop, [36](#)
- billboard, [2](#)
- chop, [3](#)
- coalesce, [27](#)
- complete, [4](#)
- complete(), [7](#)
- construction, [5](#)
- crossing (expand), [7](#)
- dplyr::count(), [33](#)
- dplyr::group\_by(), [18](#)
- dplyr::left\_join(), [4](#)
- dplyr::select(), [6, 10, 12, 20, 30, 34](#)
- drop\_na, [6](#)
- expand, [7](#)
- expand(), [4](#)
- expand.grid, [9](#)
- expand\_grid, [8](#)
- expand\_grid(), [7](#)
- extract, [9](#)
- extract(), [23, 29](#)
- fill, [10](#)
- fish\_encounters, [11](#)
- full\_seq, [11](#)
- gather, [12](#)
- gather(), [24](#)
- hoist, [13](#)
- na\_if, [27](#)
- nest, [17](#)
- nest(), [3, 20](#)
- nest\_legacy, [19](#)
- nest\_legacy(), [18](#)
- nesting (expand), [7](#)
- pack, [21](#)
- pivot\_longer, [22](#)
- pivot\_longer(), [24](#)
- pivot\_wider, [24](#)
- pivot\_wider(), [22](#)
- pivot\_wider\_spec(), [25](#)
- population (who), [35](#)
- purrr::pluck(), [14](#)
- quasiquote, [10, 12, 20, 28, 31, 34](#)
- recode, [27](#)
- relig\_income, [26](#)
- replace\_na, [27](#)
- replace\_na(), [4](#)
- rlang::as\_function(), [15, 18, 22](#)
- rlang::ensym(), [12, 20, 34](#)
- separate, [28](#)
- separate(), [10, 23, 34](#)
- separate\_rows, [29](#)
- smiths, [30](#)
- spread, [31](#)
- spread(), [25](#)
- table1, [32](#)
- table2 (table1), [32](#)
- table3 (table1), [32](#)
- table4a (table1), [32](#)
- table4b (table1), [32](#)
- table5 (table1), [32](#)
- tidyr\_legacy, [15, 18, 22](#)
- tidyselect::vars\_pull(), [10, 28, 31](#)
- tidyselect::vars\_select(), [6, 13, 30, 34](#)
- type.convert(), [10, 12, 28, 30, 31](#)
- unchop (chop), [3](#)
- uncount, [33](#)
- unite, [33](#)
- unite(), [29](#)

`unnest (nest)`, [17](#)  
`unnest()`, [20](#)  
`unnest_auto (hoist)`, [13](#)  
`unnest_legacy (nest_legacy)`, [19](#)  
`unnest_legacy()`, [18](#)  
`unnest_longer (hoist)`, [13](#)  
`unnest_wider (hoist)`, [13](#)  
`unpack (pack)`, [21](#)  
`us_rent_income`, [35](#)

`vctrs::vec_as_names()`, [15](#), [18](#), [22](#), [23](#), [25](#)

`who`, [35](#)  
`world_bank_pop`, [36](#)