

# Data manipulation basics



# Overview

## 1. The **tidyverse**

- R Packages
- Importing data

## 2. The **dplyr** package

- `filter()`
- `mutate()`
- `ifelse()`
- pipes `|>`
- `summarize()`
- `group_by()`

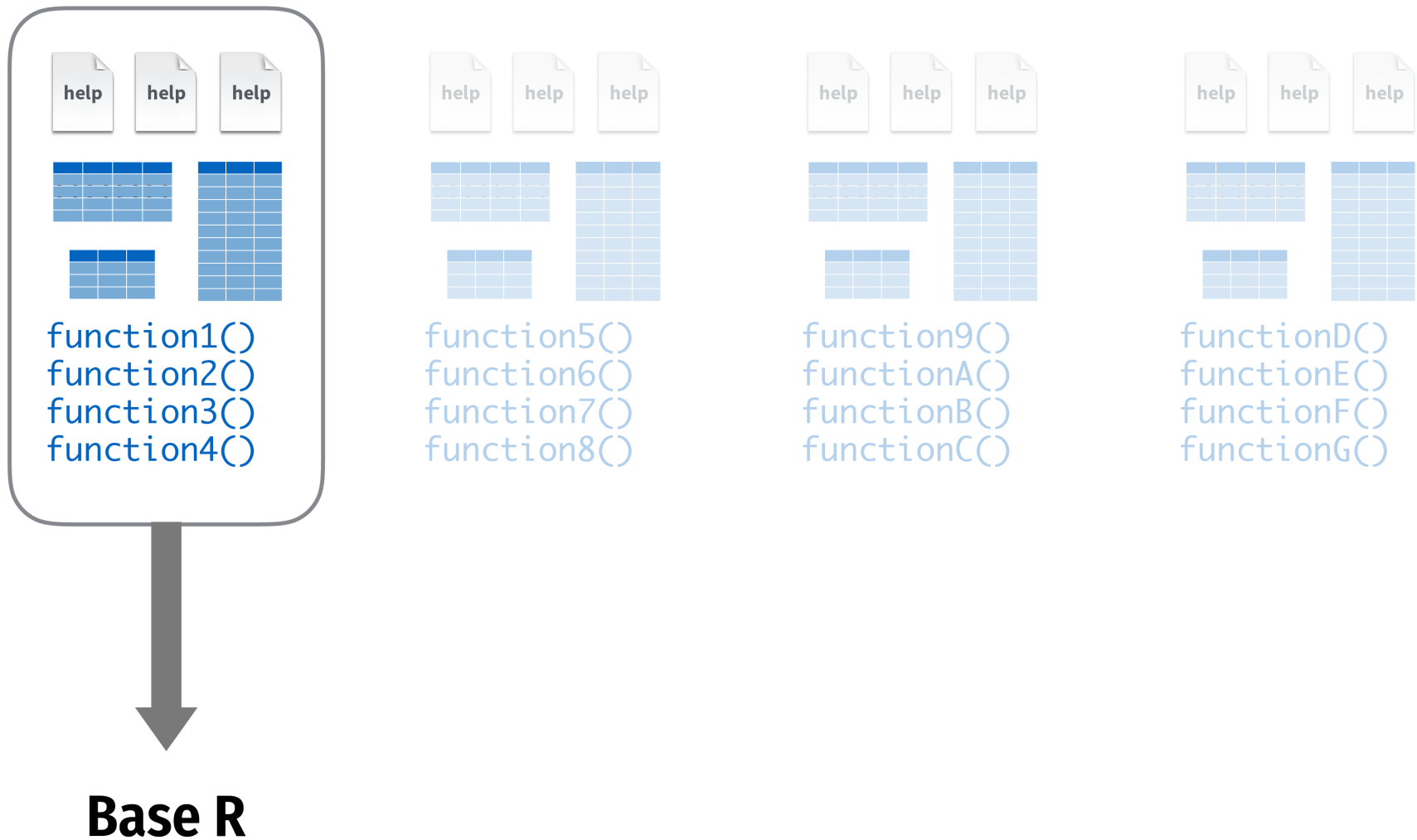
## 3. The **tidy** data format



# The tidyverse



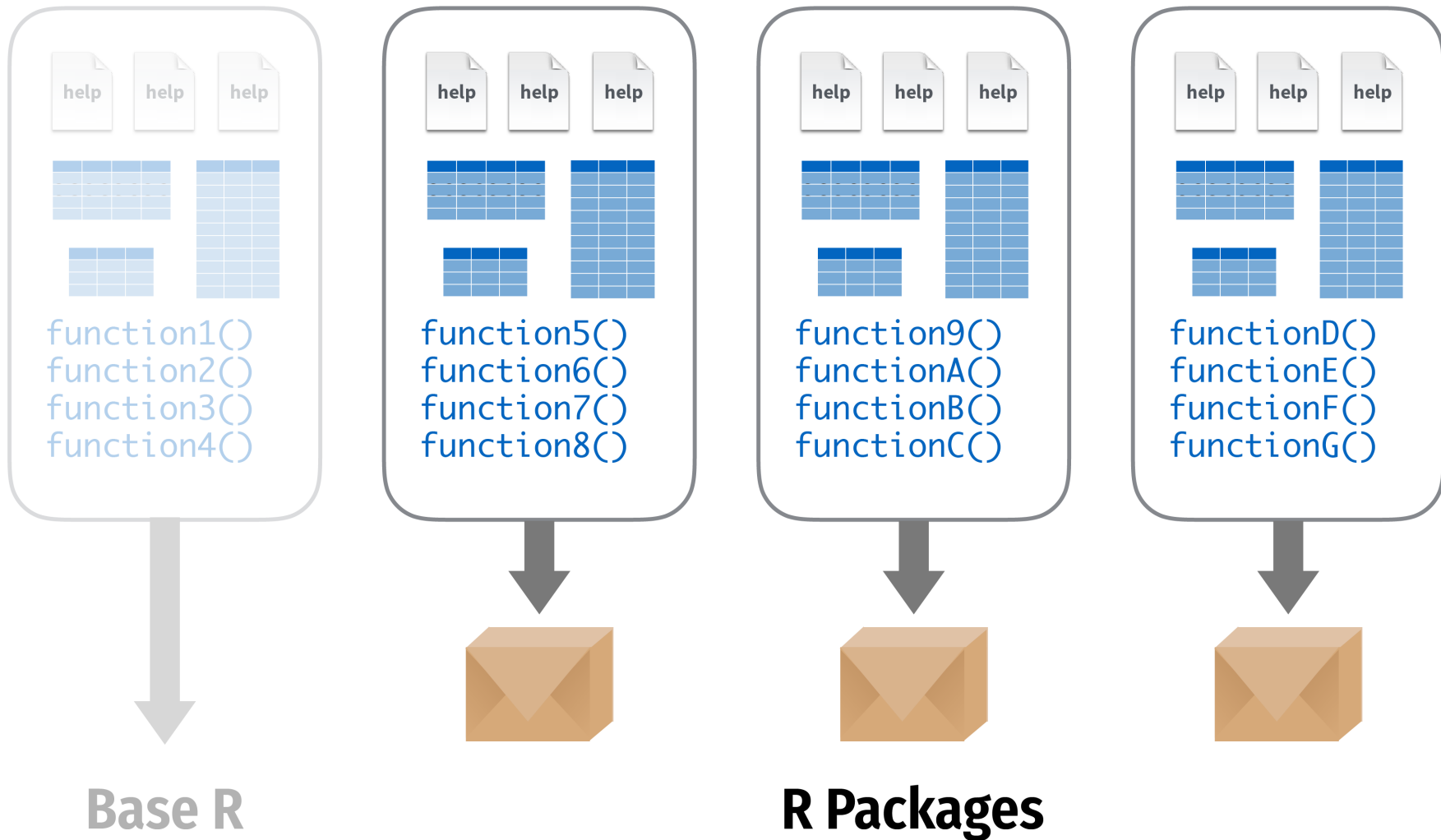
# Packages







# Packages





# Packages

- So far we only used functions that are directly available in R
  - But there are tons of user-created functions out there that can make your life so much easier
  - These functions are shared in what we call packages
- Packages are bundles of functions that R users put at the disposal of other R users
  - Packages are centralized on the Comprehensive R Archive Network (CRAN)
  - To download and install a CRAN package you can simply type `install.packages()`



# Using packages

```
1 install.packages("name")
```

- files to your computer
- Do this once per computer

```
1 library("name")
```

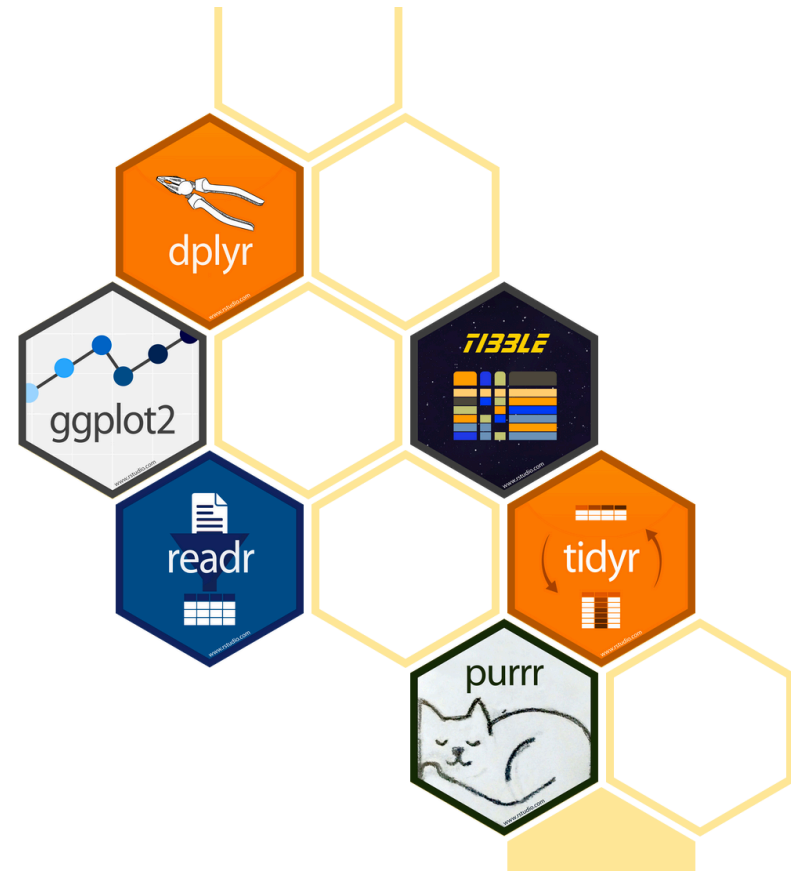
- Loads the package
- Do this once per R session



# The tidyverse

“The tidyverse is an opinionated collection of R packages designed for data science. All packages share an underlying design philosophy, grammar, and data structures.”

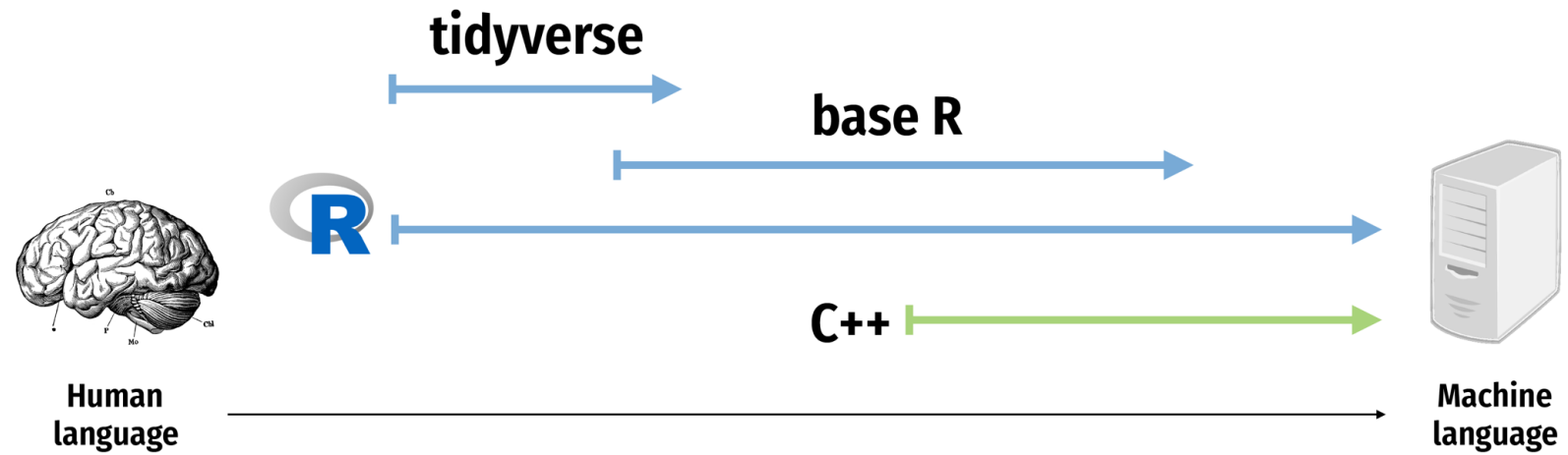
... the tidyverse makes data science faster, easier and more fun...







# The tidyverse





# The tidyverse

```
1 library(tidyverse)
```

The tidyverse package is a shortcut for installing and loading all the key tidyverse packages



# The tidyverse

```
1 install.packages("tidyverse")
```

Installs all of these:

```
1 install.packages("ggplot2")
2 install.packages("dplyr")
3 install.packages("tidyr")
4 install.packages("readr")
5 install.packages("purrr")
6 install.packages("tibble")
7 install.packages("stringr")
8 install.packages("forcats")
9 install.packages("lubridate")
10 install.packages("hms")
11 install.packages("DBI")
12 install.packages("haven")
13 install.packages("httr")
14 install.packages("jsonlite")
15 install.packages("readxl")
16 install.packages("rvest")
17 install.packages("xml2")
18 install.packages("modelr")
19 install.packages("broom")
```

```
1 library(tidyverse)
```

Loads all of these:

```
1 library(ggplot2)
2 library(dplyr)
3 library(tidyr)
4 library(readr)
5 library(purrr)
6 library(tibble)
7 library(stringr)
8 library(forcats)
9 library(lubridate)
```



# Importing data



Work with plain  
text data

```
my_data <-  
read_csv("file.csv")
```

readr

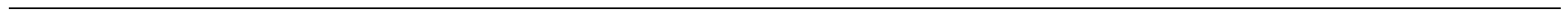
---



readxl

Work with Excel  
files

```
my_data <-  
read_excel("file.xlsx")
```







Work with Stata,  
SPSS, and SAS data

```
my_data <-  
read_stata("file.dta")
```

haven



## Data from R-Packages

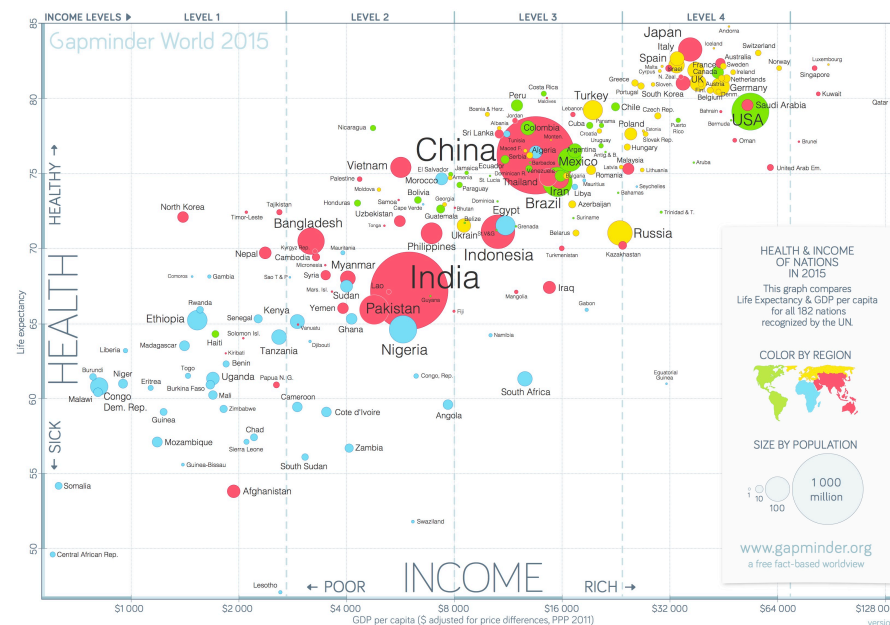
Some data sets can be downloaded as packages in R. For example, the `gapminder` data set.

## Install the package

```
1 install.packages(gapminder)
```

## Then load the data

```
1 library(gapminder)
2
3 # The data() function in R is used to list, load,
4 # and access built-in or package-provided datasets
5 data(gapminder)
```

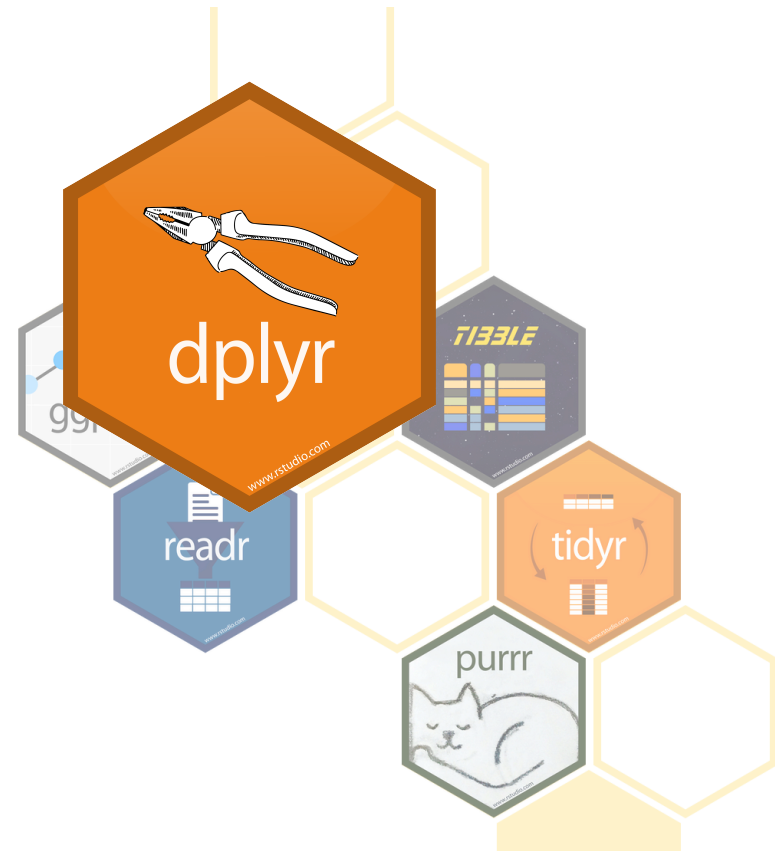




# The `dplyr` package



tidyverse

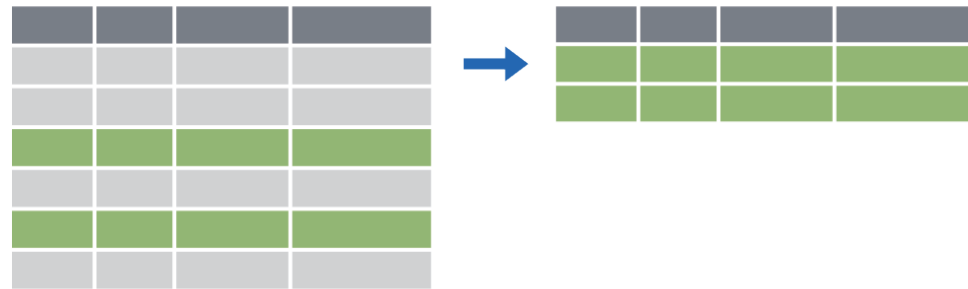






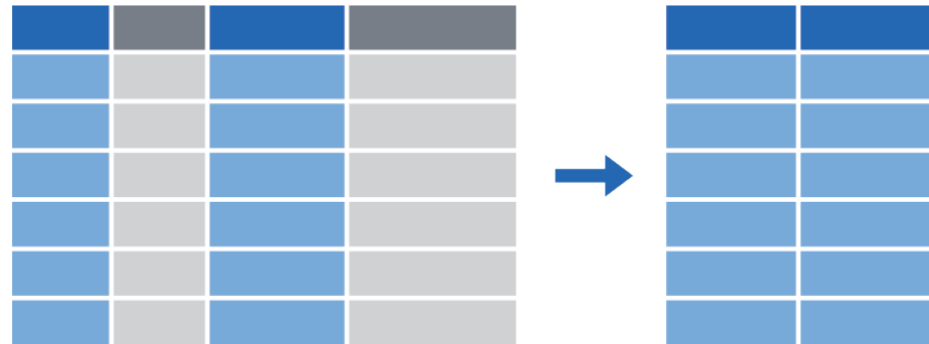
# dplyr: verbs for manipulating data

Extract rows with  
`filter()`



filter

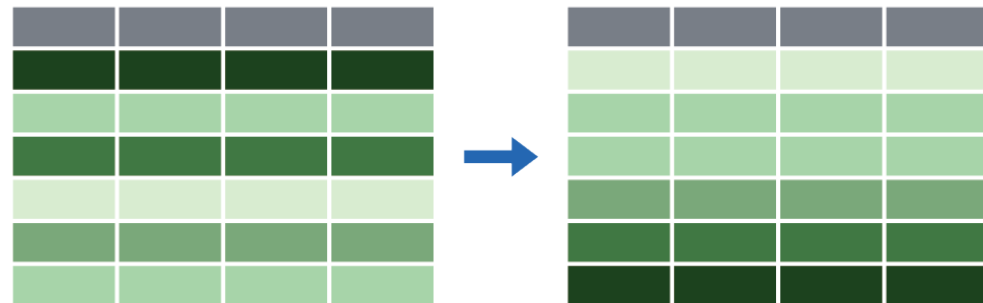
---



Extract columns with  
`select()`

select

---

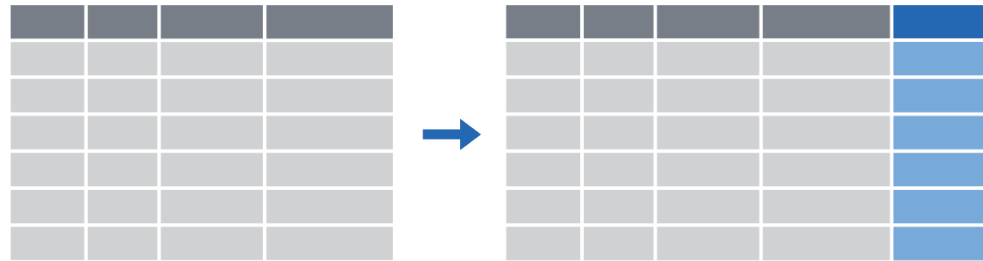


Arrange/sort rows with  
`arrange()`

arrange

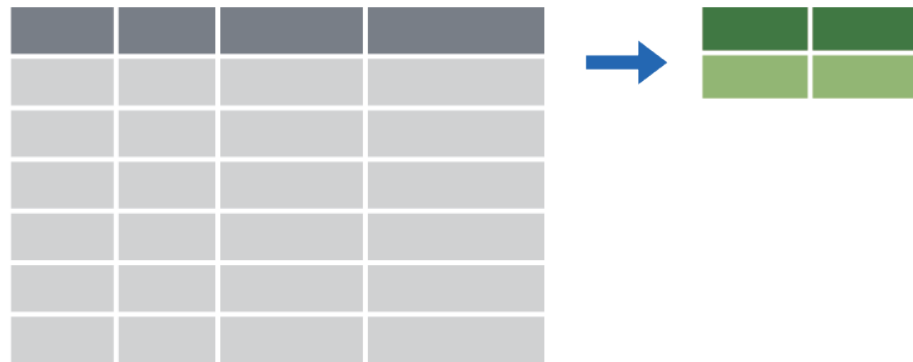
---

Make new columns with  
`mutate()`



`mutate`

Make group summaries  
with  
`group_by() |>`  
`summarize()`



`summarize`



# filter()

Extract rows that meet some sort of test

The general idea:

```
1 filter(  
2   some_data,  
3   ... # one or more tests  
4 )
```

Let's try this on the gapminder data set that you've installed earlier.

```
1 filter(.data = gapminder, country == "Denmark")
```

country	continent	year
Denmark	Europe	1952
Denmark	Europe	1957
Denmark	Europe	1962
Denmark	Europe	1967
Denmark	Europe	1972
...	...	...



# Logical tests

Test	Meaning	Test	Meaning
<code>x &lt; y</code>	Less than	<code>x %in% y</code>	In (group membership)
<code>x &gt; y</code>	Greater than	<code>is.na(x)</code>	Is missing
<code>==</code>	Equal to	<code>!is.na(x)</code>	Is not missing
<code>x &lt;= y</code>	Less than or equal to		
<code>x &gt;= y</code>	Greater than or equal to		
<code>x != y</code>	Not equal to		





# Your turn #1: Filtering

04:00

Use `filter()` and logical tests to show...

1. The data for Canada
2. All data for countries in Oceania
3. Rows where the life expectancy is greater than 82



# Common Mistakes

Using `=` instead of `==`

Bad

```
1 filter(gapminder, country = "Canada")
```

Good

```
1 filter(gapminder, country == "Canada")
```

Forgetting quotes (`""`)

Bad

```
1 filter(gapminder, country == Canada)
```

Good

```
1 filter(gapminder, country == "Canada")
```



# filter() with multiple conditions

Extract rows that meet *every* test

```
1 filter(gapminder, country == "Denmark", year > 2000)
```

# A tibble: 2 × 6

	country	continent	year	lifeExp	pop	gdpPercap
	<fct>	<fct>	<int>	<dbl>	<int>	<dbl>
1	Denmark	Europe	2002	77.2	5374693	32167.
2	Denmark	Europe	2007	78.3	5468120	35278.



# Boolean operators

Operator	Meaning
<code>a &amp; b</code>	and
<code>a   b</code>	or
<code>!a</code>	not





# Boolean operators

The default is “and”

These do the same thing:

```
1 filter(gapminder,  
2       country == "Denmark",  
3       year > 2000)
```

```
1 filter(gapminder,  
2       country == "Denmark" &  
3       year > 2000)
```



# Your turn #2: Filtering

04:00

Use `filter()` and Boolean logical tests to show...

1. Canada before 1970
2. Countries where life expectancy in 2007 is below 50
3. Countries where life expectancy in 2007 is below 50 and are not in Africa



# Common Mistakes

Collapsing multiple tests into one

Bad

```
1 filter(gapminder,  
2       1960 < year < 1980)
```

Good

```
1 filter(gapminder,  
2       year > 1960,  
3       year < 1980)
```

Using multiple tests instead of `%in%`

Bad

```
1 filter(gapminder,  
2       country == "Mexico",  
3       country == "Canada",  
4       country == "United States")
```

Good

```
1 filter(gapminder,  
2       country %in% c("Mexico", "Canada",  
3                     "United States"))
```



# Common Syntax

Every `dplyr` verb function follows the same pattern

```
1 verb(data, ...)
```

`verb` = dplyr function/verb

`data` = data frame to transform

`...` = what you tell the verb to do exactly





# mutate()

Create new columns

The general idea:

```
1 mutate(  
2   some_data,  
3   ... # new columns to make  
4 )
```

Let's try this on the gapminder data

```
1 mutate(gapminder, gdp = gdpPercap * pop)
```

country	year	...	gdp
Afghanistan	1952	...	6567086330
Afghanistan	1957	...	7585448670
Afghanistan	1962	...	8758855797
Afghanistan	1967	...	9648014150
Afghanistan	1972	...	9678553274
Afghanistan	1977	...	11697659231



# mutate()

Create new columns

The general idea:

```
1 mutate(  
2   some_data,  
3   ... # new columns to make  
4 )
```

We can also create multiple new columns at once

```
1 mutate(gapminder, gdp = gdpPercap * pop,  
2           pop_mil = round(pop / 1000000))
```

country	year	...	gdp
Afghanistan	1952	...	6567086330
Afghanistan	1957	...	7585448670
Afghanistan	1962	...	8758855797
Afghanistan	1967	...	9648014150
Afghanistan	1972	...	9678553274
Afghanistan	1977	...	11697659231



# ifelse()

Do conditional tests within `mutate()`

```
1 ifelse(test,  
2       value_if_true,  
3       value_if_false)
```

`test` = a logical test

`value_if_true` = what happens if test is true

`value_if_false` = what happens if test is false



# ifelse()

The new variable can take any sort of class

```
1 # a new logical variable
2 mutate(gapminder,
3         after_1960 = ifelse(year > 1960, TRUE, FALSE)
4         )
```

```
1 # a new character variable
2 mutate(gapminder,
3         after_1960 = ifelse(year > 1960,
4                             "After 1960",
5                             "Before 1960")
6         )
```

```
1 # a new numeric variable
2 mutate(gapminder,
3         after_1960 = ifelse(year > 1960, 0, 1)
4         )
```





# Your turn #3: Mutating

05:00

Use `mutate()` to...

1. Add an `africa` column that is TRUE if the country is on the African continent
2. Add a column for logged GDP per capita (hint: use `log()`)
3. Add an `africa_asia` column that says “Africa or Asia” if the country is in Africa or Asia, and “Not Africa or Asia” if it’s not



# What if you have multiple verbs?

## Solution 1: Intermediate variables

```
1 gapminder_2002 <- filter(gapminder, year == 2002)
2
3 gapminder_2002_log <- mutate(gapminder_2002,
4                               log_gdpPercap = log(gdpPercap))
```

## Solution 2: Nested functions

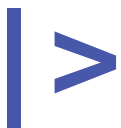
```
1 filter(mutate(gapminder_2002,
2               log_gdpPercap = log(gdpPercap)),
3        year == 2002)
```

## Solution 3: Pipes!

- The `|>` operator (pipe) takes an object on the left and passes it as the first argument of the function on the right

```
1 gapminder |>
2   filter(year == 2002) |>
3   mutate(log_gdpPercap = log(gdpPercap))
```





## Why using pipes?

```
1 leave_house(get_dressed(get_out_of_bed(wake_up(me, time = "8:00"), side = "correct"),
2           pants = TRUE, shirt = TRUE), car = TRUE, bike = FALSE)
```

... 🤪 not easy to read

```
1 me |>
2   wake_up(time = "8:00") |>
3   get_out_of_bed(side = "correct") |>
4   get_dressed(pants = TRUE, shirt = TRUE) |>
5   leave_house(car = TRUE, bike = FALSE)
```

... 🎉 easy to read



# |> vs %>%

- There are actually multiple pipes!
- %>% was invented first, but requires a package to use
- |> is part of base R
- They're interchangeable 99% of the time (Just be consistent)



You do not have to type the pipe by hand every time

You can use the shortcut `cmd + shift + m` in R Studio.





# summarize()

Compute a table of summaries

## 1. Take a data frame

country	continent	year	lifeExp
Afghanistan	Asia	1952	28.801
Afghanistan	Asia	1957	30.332
Afghanistan	Asia	1962	31.997
Afghanistan	Asia	1967	34.02
...	...	...	...

## 2. Make a summary

```
1 gapminder |>
2   summarize(mean_life = mean(lifeExp))
```

```
# A tibble: 1 × 1
  mean_life
    <dbl>
1    59.5
```

## Or several summaries

```
1 gapminder |>
2   summarize(mean_life = mean(lifeExp),
3             min_life = min(lifeExp))
```

```
# A tibble: 1 × 2
  mean_life min_life
    <dbl>    <dbl>
1    59.5     23.6
```



# Your turn #4: Summarizing

05:00

Use `summarize()` to calculate...

1. The first (minimum) year in the dataset
2. The last (maximum) year in the dataset
3. The number of rows in the dataset (use the cheatsheet)
4. The number of distinct countries in the dataset (use the cheatsheet)



# Your turn #5: Summarizing

05:00

Use `filter()` and `summarize()` to calculate...

1. the number of unique countries and
2. the median life expectancy

on the African continent in 2007.



# group\_by()

Put rows into groups based on values in a column

```
1 gapminder |> group_by(continent)
```

- Nothing happens by itself!
- Powerful when combined with `summarize()`





# group\_by()

country	continent	year	lifeExp
Afghanistan	Asia	1952	28.801
Afghanistan	Asia	1957	30.332
Afghanistan	Asia	1962	31.997
Afghanistan	Asia	1967	34.02
...	...	...	...

## A simple summary

```
1 gapminder |>
2   summarize(n_countries = n_distinct(country))
```

```
# A tibble: 1 × 1
  n_countries
    <int>
1       142
```

## A grouped summary

```
1 gapminder |>
2   group_by(continent) |>
3   summarize(n_countries = n_distinct(country))
```

```
# A tibble: 5 × 2
  continent n_countries
  <fct>      <int>
1 Africa         52
2 Americas        25
3 Asia           33
4 Europe         30
5 Oceania         2
```



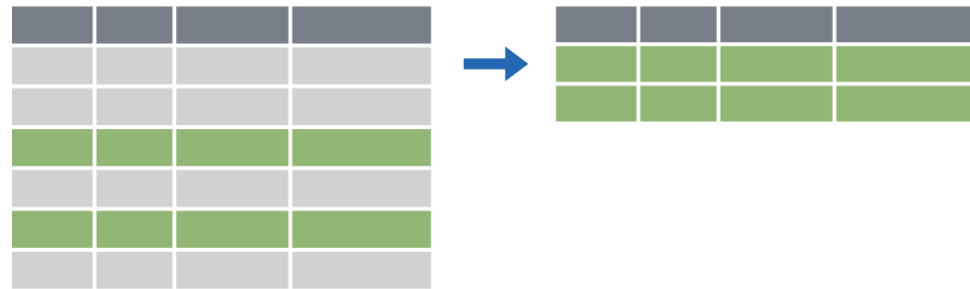
# Your turn #6: Grouping and summarizing

1. Find the minimum, maximum, and median life expectancy for each continent
2. Find the minimum, maximum, and median life expectancy for each continent in 2007 only



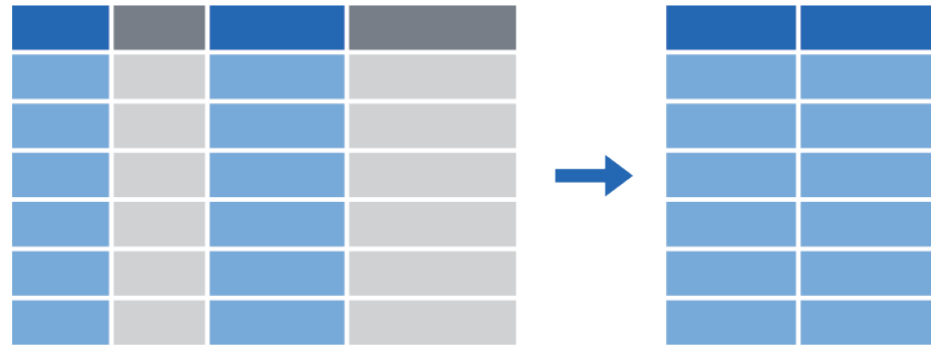
# dplyr: verbs for manipulating data

Extract rows with  
`filter()`



filter

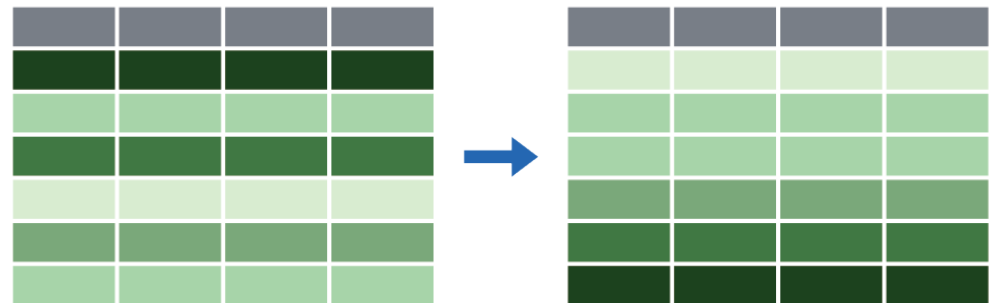
---



Extract columns with  
`select()`

select

---

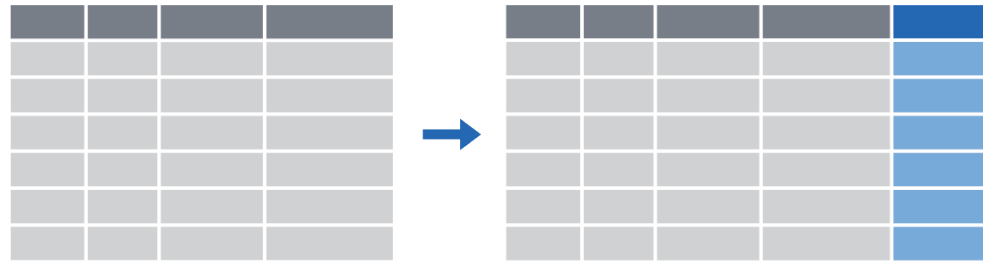


Arrange/sort rows with  
`arrange()`

arrange

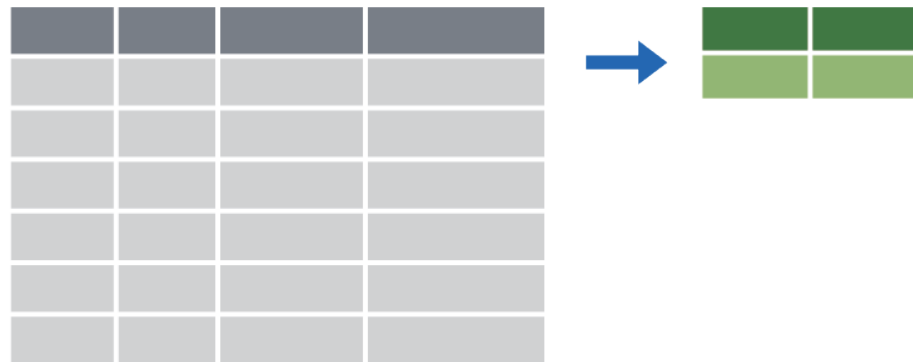
---

Make new columns with  
`mutate()`



`mutate`

Make group summaries  
with  
`group_by() |>`  
`summarize()`



`summarize`





# Tidy data



# Tidy data

You can represent the same underlying data in multiple ways.

```
# A tibble: 6 × 4
  country      year cases
population
  <chr>      <dbl> <dbl>
<dbl>
1 Afghanistan 1999     745
19987071
2 Afghanistan 2000    2666
20595360
3 Brazil      1999   37737
172006362
4 Brazil      2000   80488
174504898
5 China       1999  212258
1272915272
6 China       2000  213766
1280428583
```

```
# A tibble: 12 × 4
  country      year type
count
  <chr>      <dbl> <chr>
<dbl>
1 Afghanistan 1999 cases
745
2 Afghanistan 1999 population
19987071
3 Afghanistan 2000 cases
2666
4 Afghanistan 2000 population
20595360
5 Brazil      1999 cases
37737
6 Brazil      1999 population
172006362
7 Brazil      2000 cases
80488
8 Brazil      2000 population
174504898
9 China       1999 cases
212258
```

```
# A tibble: 6 × 3
  country      year rate
  <chr>      <dbl> <chr>
1 Afghanistan 1999 745/19987071
2 Afghanistan 2000 2666/20595360
3 Brazil      1999 37737/172006362
4 Brazil      2000 80488/174504898
5 China       1999 212258/1272915272
6 China       2000 213766/1280428583
```



# Tidy data

Tidy data has the following properties:

```
# A tibble: 6 × 4
  country    year cases
population
  <chr>      <dbl> <dbl>
<dbl>
1 Afghanistan 1999     745
19987071
2 Afghanistan 2000    2666
20595360
3 Brazil      1999   37737
172006362
4 Brazil      2000   80488
174504898
5 China       1999  212258
1272915272
6 China       2000  213766
1280428583
```

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280428583

variables

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280428583

observations

country	year	cases	population
Afghanistan	1999	745	19987071
Afghanistan	2000	2666	20595360
Brazil	1999	37737	172006362
Brazil	2000	80488	174504898
China	1999	212258	1272915272
China	2000	213766	1280428583

values

1. Variables are columns
2. Observations are rows
3. Values are cells



# Why ensure that your data is tidy?

There are two main advantages:

1. There's a general advantage to picking one consistent way of storing data. If you have a consistent data structure, it's easier to learn the tools that work with it because they have an underlying uniformity.
2. There's a specific advantage to placing variables in columns because it allows R's vectorized nature to shine. As you learned in `?@sec-mutate` and `?@sec-summarize`, most built-in R functions work with vectors of values. That makes transforming tidy data feel particularly natural.

`dplyr`, `ggplot2`, and all the other packages in the tidyverse are designed to work with tidy data.





# Will I ever encounter a dataset that isn't tidy?

Yes, unfortunately, most real data is untidy.

There are two main reasons:

1. Data is often organized to facilitate some goal other than analysis. For example, it's common for data to be structured to make data entry, not analysis, easy.
2. Most people aren't familiar with the principles of tidy data, and it's hard to derive them yourself unless you spend a lot of time working with data.



# Pivoting data

`tidyr` provides two main functions to “pivot” data in a tidy format:  
`pivot_longer()` and `pivot_wider()`

Here, we’ll only discuss `pivot_longer()` because it’s the most common case.



# pivot\_longer()

- Suppose we have three patients with `ids` A, B, and C, and we take two blood pressure measurements on each patient.
- We'll create the data with `tribble()`, a handy function for constructing small tibbles by hand:

```
1 df <- tribble(  
2   ~id, ~bp1, ~bp2,  
3   "A", 100, 120,  
4   "B", 140, 115,  
5   "C", 120, 125  
6 )  
7  
8 df
```

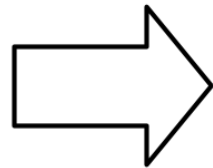
```
# A tibble: 3 × 3  
  id      bp1  bp2  
<chr> <dbl> <dbl>  
1 A      100   120  
2 B      140   115  
3 C      120   125
```



# pivot\_longer()

- We want our new dataset to have three variables: **id** (already exists), **measurement** (the column names), and **value** (the cell values)
- To achieve this, we need to pivot **df** longer

id	bp1	bp2
A	100	120
B	140	115
C	120	125



id	measurement	value
A	bp1	100
A	bp2	120
B	bp1	140
B	bp2	115
C	bp1	120
C	bp2	125

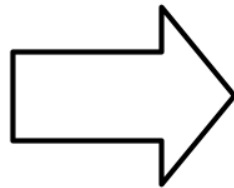




# pivot\_longer()

- The values in a column that was already a variable in the original dataset (*id*) need to be repeated, once for each column that is pivoted.

id	bp1	bp2
A	100	120
B	140	115
C	120	125



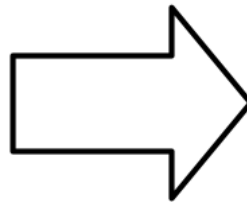
id	measurement	value
A	bp1	100
A	bp2	120
B	bp1	140
B	bp2	115
C	bp1	120
C	bp2	125



# pivot\_longer()

- The column names become values of the new variable `measurement`

id	bp1	bp2
A	100	120
B	140	115
C	120	125



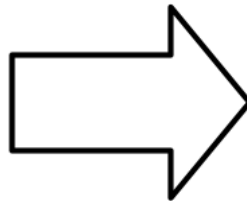
id	measurement	value
A	bp1	100
A	bp2	120
B	bp1	140
B	bp2	115
C	bp1	120
C	bp2	125



# pivot\_longer()

- The cell values become values of the new variable `value`

id	bp1	bp2
A	100	120
B	140	115
C	120	125



id	measurement	value
A	bp1	100
A	bp2	120
B	bp1	140
B	bp2	115
C	bp1	120
C	bp2	125



# pivot\_longer()

```
# A tibble: 3 × 3
  id      bp1    bp2
<chr> <dbl> <dbl>
1 A      100    120
2 B      140    115
3 C      120    125
```

```
1 df |>
2   pivot_longer(
3     cols = bp1:bp2,
4     names_to = "measurement",
5     values_to = "value"
6   )
```

```
# A tibble: 6 × 3
  id      measurement value
<chr> <chr>         <dbl>
1 A      bp1             100
2 A      bp2             120
3 B      bp1             140
4 B      bp2             115
5 C      bp1             120
6 C      bp2             125
```

After the data, there are three key arguments:

- `cols` specifies which columns need to be pivoted, i.e. which columns aren't variables. This argument uses the same syntax as `select()`
- `names_to` names the variable in which column names should be stored
- `values_to` names the variable in which cell values should be stored





# Your turn #7: Pivoting

The `billboard` dataset which comes with the `tidyverse` package records the billboard rank of songs in the year 2000.

```
1 head(billboard)
```

```
# A tibble: 6 × 79
  artist      track date.entered  wk1  wk2  wk3  wk4  wk5  wk6  wk7  wk8
  <chr>      <chr> <date>      <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
1 2 Pac      Baby... 2000-02-26    87   82   72   77   87   94   99   NA
2 2Ge+her    The ... 2000-09-02    91   87   92   NA   NA   NA   NA   NA
3 3 Doors Do... Kryp... 2000-04-08    81   70   68   67   66   57   54   53
4 3 Doors Do... Loser 2000-10-21    76   76   72   69   67   65   55   59
5 504 Boyz    Wobb... 2000-04-15    57   34   25   17   17   31   36   49
6 98^0       Give... 2000-08-19    51   39   34   26   26   19    2    2

# i 68 more variables: wk9 <dbl>, wk10 <dbl>, wk11 <dbl>, wk12 <dbl>,
# wk13 <dbl>, wk14 <dbl>, wk15 <dbl>, wk16 <dbl>, wk17 <dbl>, wk18 <dbl>,
# wk19 <dbl>, wk20 <dbl>, wk21 <dbl>, wk22 <dbl>, wk23 <dbl>, wk24 <dbl>,
# wk25 <dbl>, wk26 <dbl>, wk27 <dbl>, wk28 <dbl>, wk29 <dbl>, wk30 <dbl>,
# wk31 <dbl>, wk32 <dbl>, wk33 <dbl>, wk34 <dbl>, wk35 <dbl>, wk36 <dbl>,
# wk37 <dbl>, wk38 <dbl>, wk39 <dbl>, wk40 <dbl>, wk41 <dbl>, wk42 <dbl>,
# wk43 <dbl>, wk44 <dbl>, wk45 <dbl>, wk46 <dbl>, wk47 <dbl>, wk48 <dbl>, ...
```



# Your turn #7: Pivoting

1. Use `pivot_longer()` to tidy the data (Tip: Create the new variables `week` and `rank`). Assign the resulting data frame to a new data frame called `tidy_billboard`.
2. Use the new `tidy_billboard` data frame to calculate which song has been the longest on rank 1 (Tip: use `filter()`, `group_by()` and `summarize()`)



**That's it for today :)**

