

Data wrangling with dplyr

Week 4 (10/22/25)



Artwork by Allison Horst

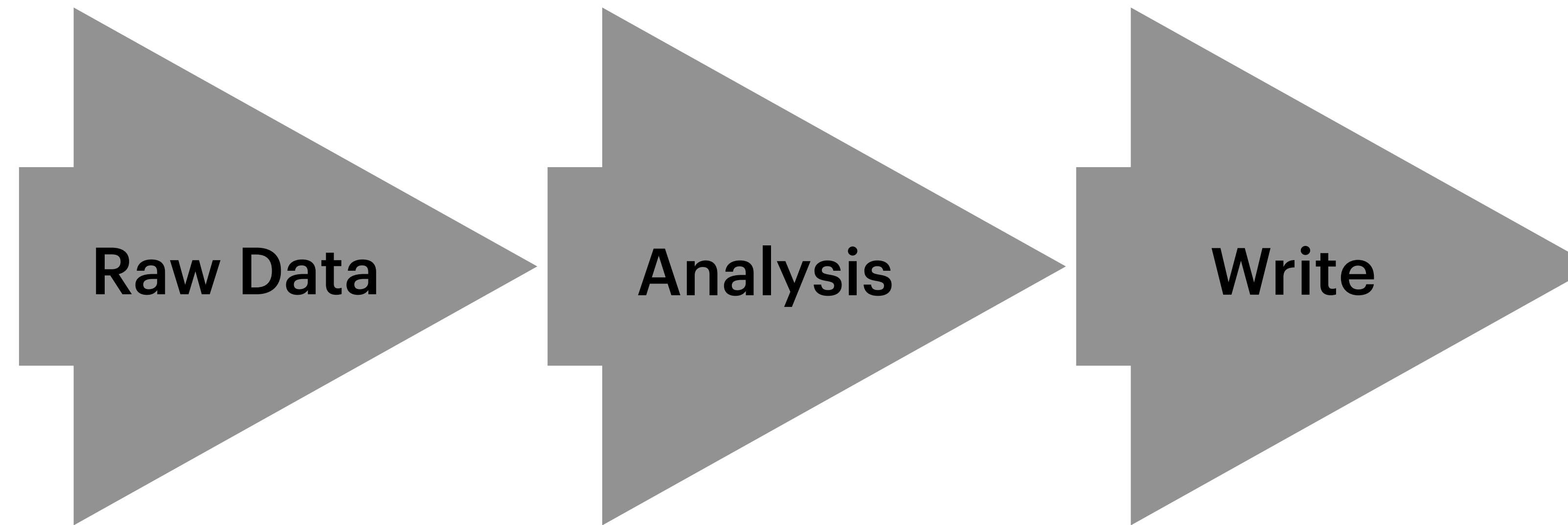
Stepfanie M. Aguillon

Outline of today's class

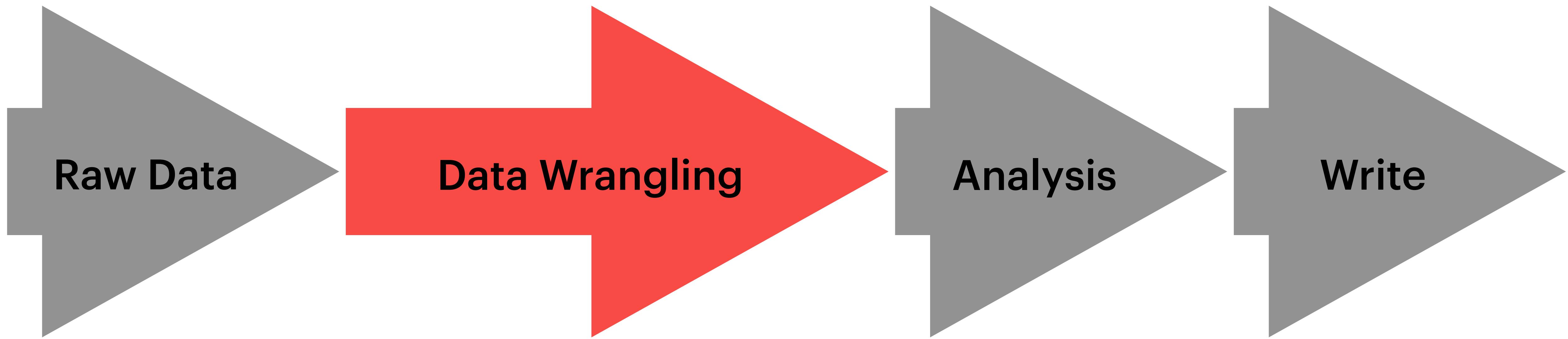
- What is data wrangling?
- Introduction to dplyr
- Continue practicing with RStudio/GitHub and Quarto

What is data wrangling?

The research pipeline

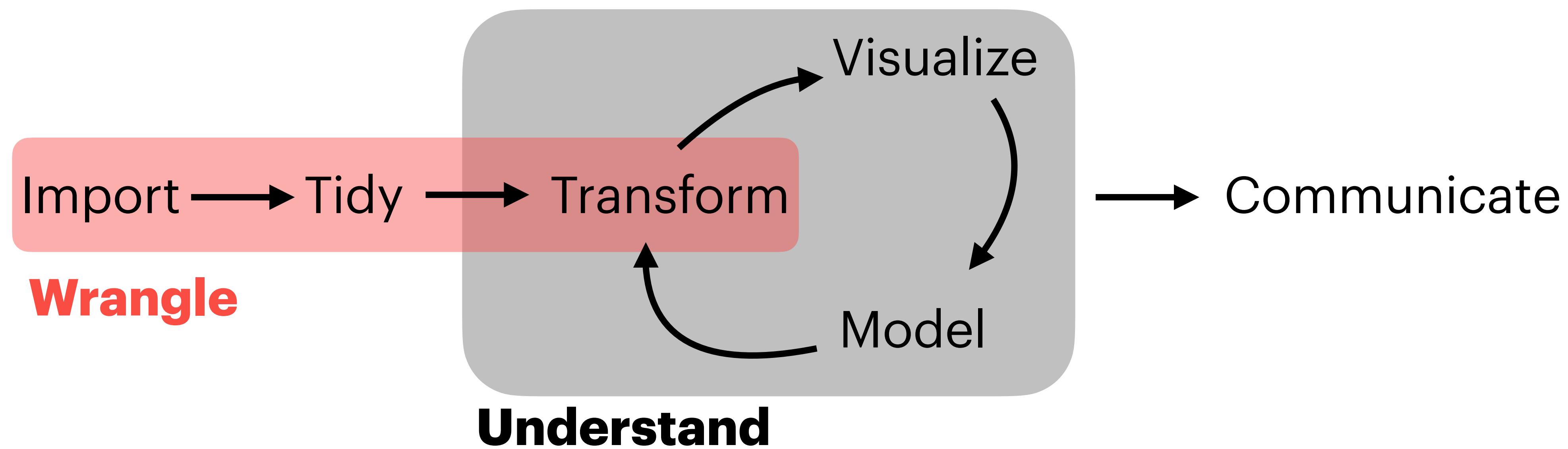


The research pipeline



Data wrangling = Getting your data into a useable form for downstream analyses and data visualization

Data wrangling



Introduction to dplyr

tidyverse package
for manipulating data

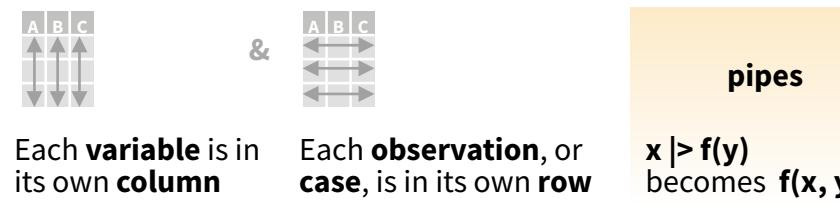


dplyr Cheatsheet

Data transformation with dplyr :: CHEATSHEET

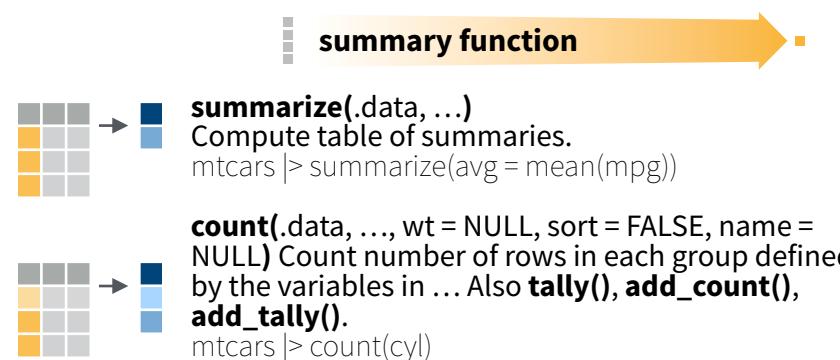


dplyr functions work with pipes and expect **tidy data**. In tidy data:



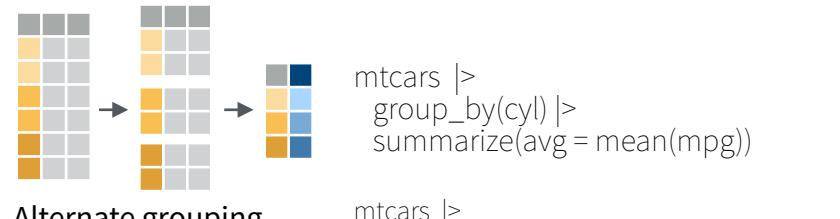
Summarize Cases

Apply **summary functions** to columns to create a new table of summary statistics. Summary functions take vectors as input and return one value (see back).



Group Cases

Use `group_by(.data, ..., .add = FALSE, .drop = TRUE)` to create a "grouped" copy of a table grouped by columns in ... dplyr functions will manipulate each "group" separately and combine the results.



Alternate grouping syntax with `.by` as an argument:
Use `rowwise(.data, ...)` to group data into individual rows. dplyr functions will compute results for each row. Also apply functions to list-columns. See tidyverse cheat sheet for list-column workflow.



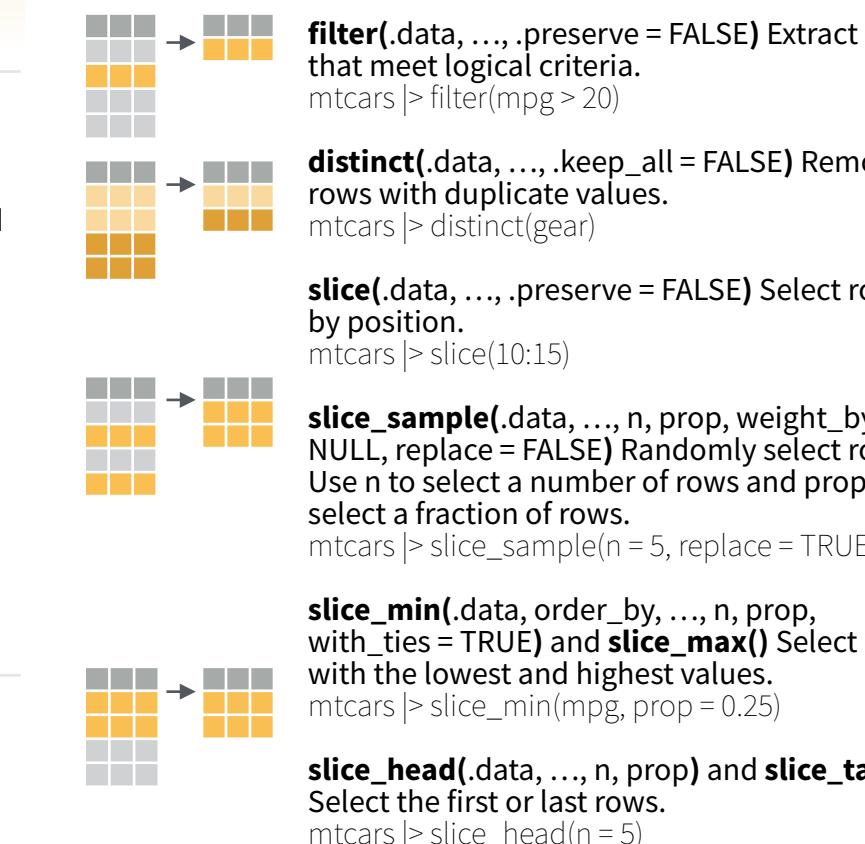
`ungroup(x, ...)` Returns ungrouped copy of table.
`g_mtcars <- mtcars |> group_by(cyl)`
`ungroup(g_mtcars)`



Manipulate Cases

EXTRACT CASES

Row functions return a subset of rows as a new table.



Logical and boolean operators to use with filter()

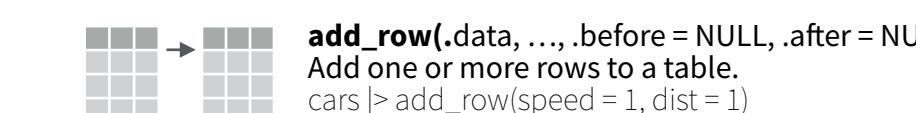
`==` `<` `<=` `is.na()` `%in%` `|` `xor()`
`!=` `>` `>=` `!is.na()` `!` `&`

See `?base::Logic` and `?Comparison` for help.

ARRANGE CASES



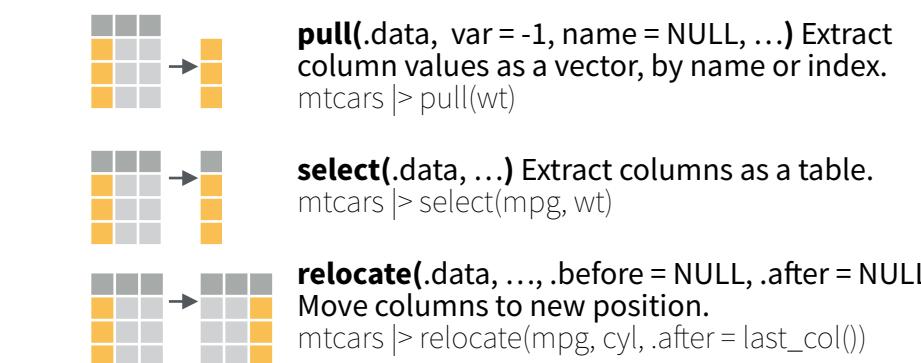
ADD CASES



Manipulate Variables

EXTRACT VARIABLES

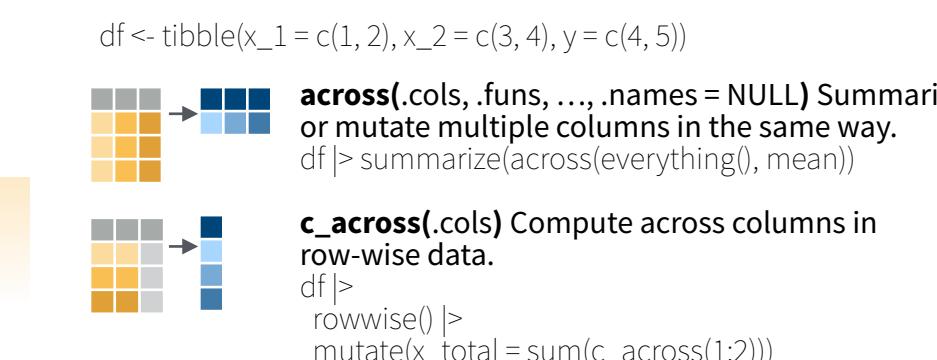
Column functions return a set of columns as a new vector or table.



Use these helpers with select() and across()

e.g. mtcars |> select(mpg:cyl)
`contains(match)` `num_range(prefix, range)` ; e.g., mpg:cyl
`ends_with(match)` `all_of(x)/any_of(x, ..., vars)` !; e.g., !gear
`starts_with(match)` `matches(match)` `everything()`

MANIPULATE MULTIPLE VARIABLES AT ONCE



MAKE NEW VARIABLES

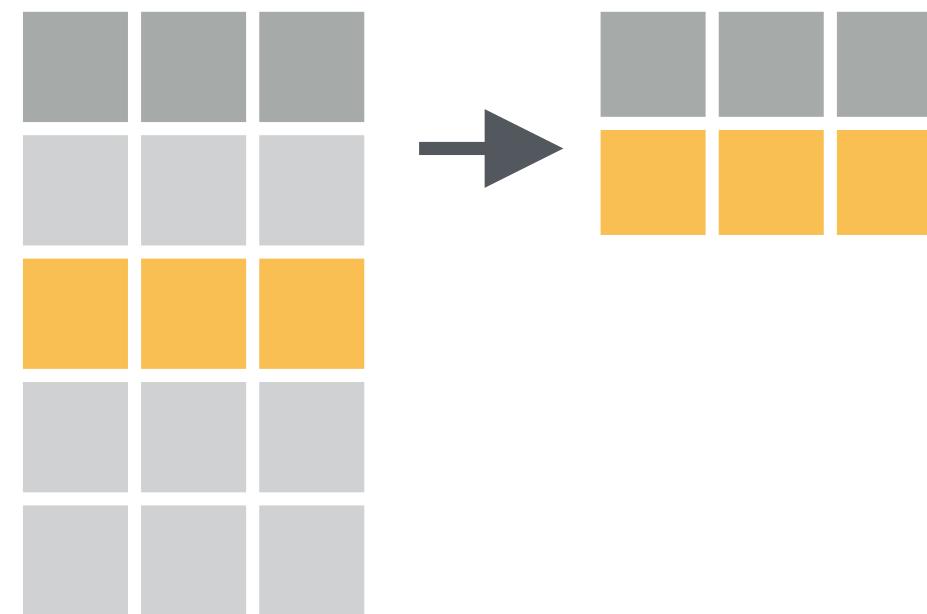
Apply **vectorized functions** to columns. Vectorized functions take vectors as input and return vectors of the same length as output (see back).



Four important dplyr functions

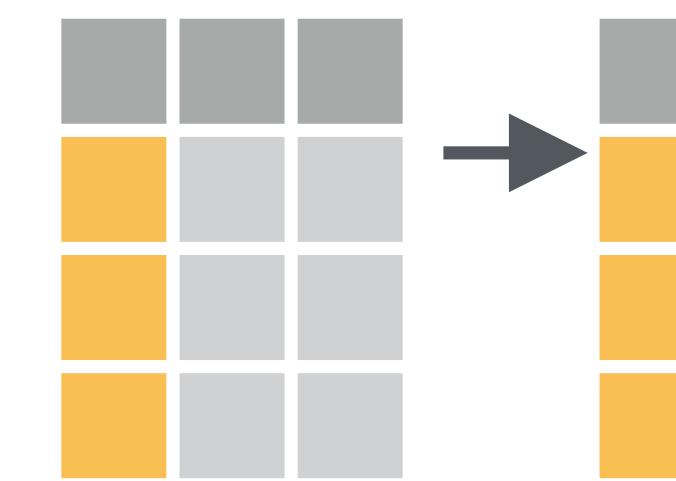
1

`filter()`



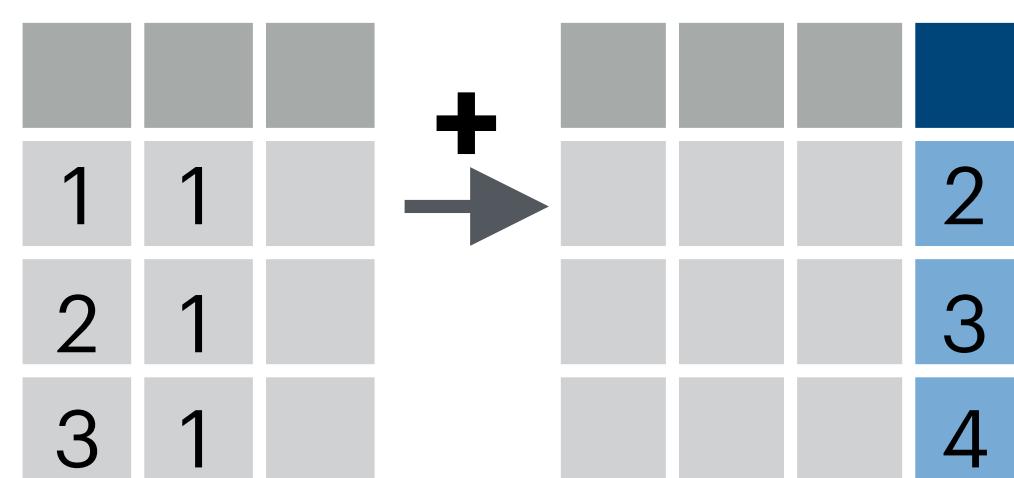
2

`select()`



3

`mutate()`



4

`arrange()`



Four important dplyr functions

1

`filter()`

filter observations
by their values

2

`select()`

select variables by
their column names

3

`mutate()`

mutate existing variables
to create new ones

4

`arrange()`

arrange observations
by their values

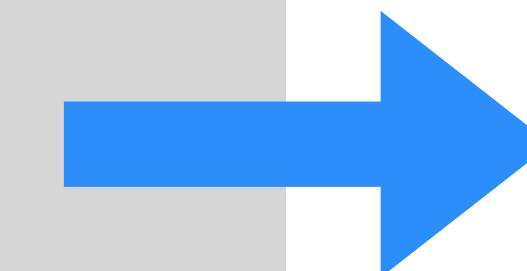
“The Pipe” is a connector

Allows you to combine a sequence of data wrangling steps

```
df_1 <- filter(df, <conditions>)
df_2 <- select(df_1, <conditions>)
df_3 <- mutate(df_2, <conditions>)
```

```
df_new <- df %>%
  filter(<conditions>) %>%
  select(<conditions>) %>%
  mutate(<conditions>)
```

the pipe!



original pipe

%>%

newer pipe

|>

You've already seen a type of connector

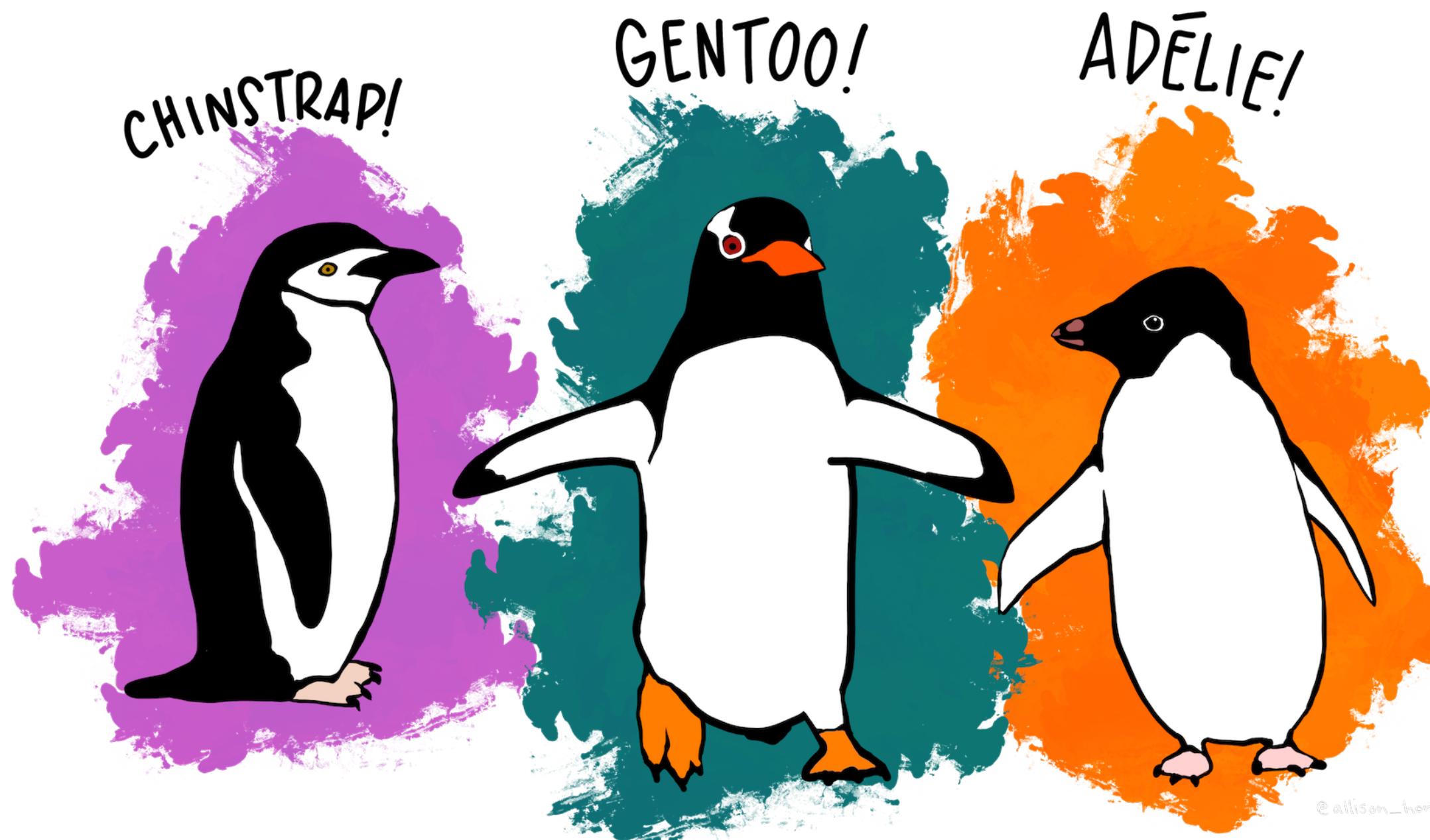
```
ggplot(data = <DATA>) +  
  geom_<TYPE>(mapping = aes(<MAPPINGS>))
```

necessary in ggplot2 to
tell R this is all one plot

Practice, practice, practice...

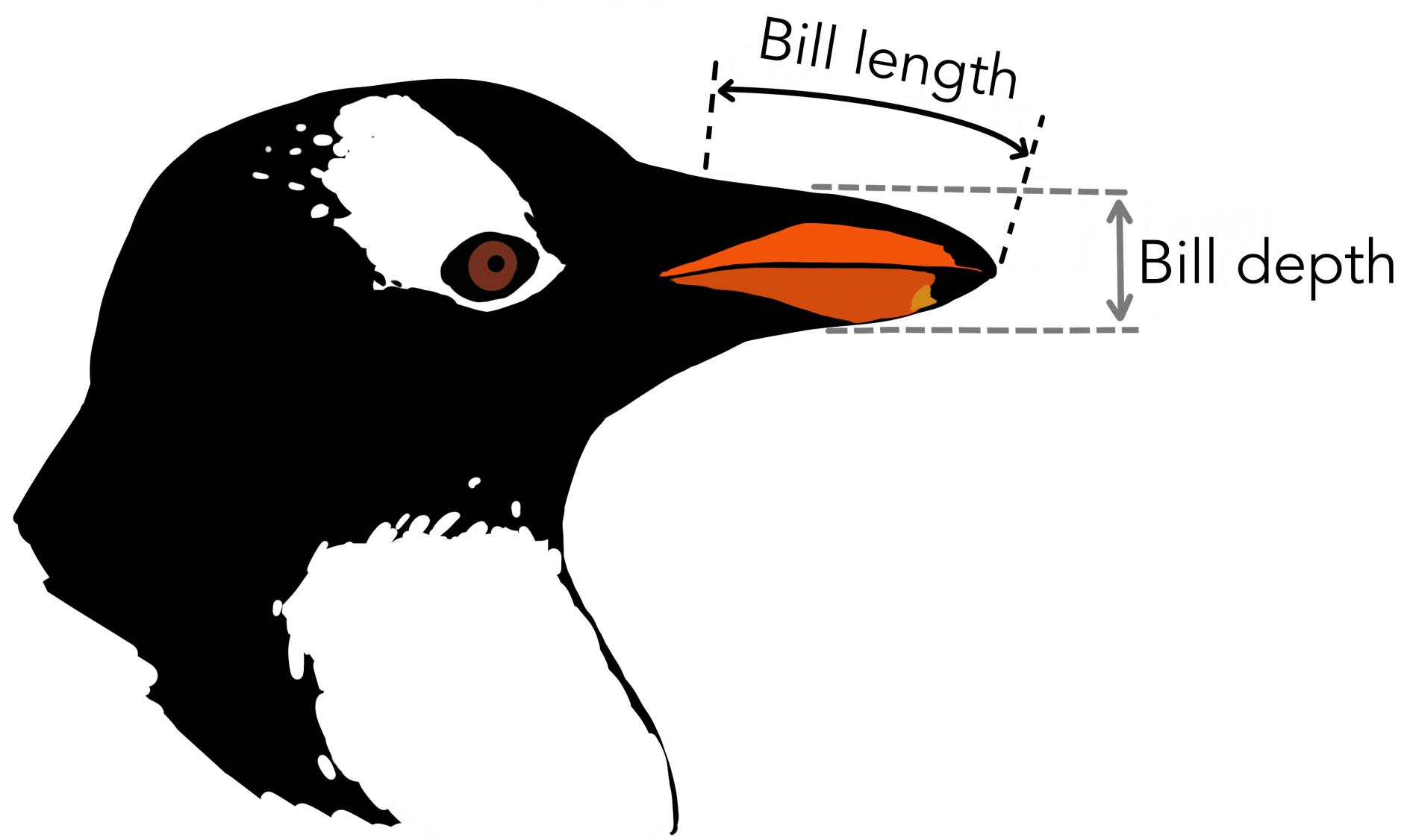
The penguins dataset

3 penguin species



Artwork by Allison Horst
© allison_horst

Morphological measures

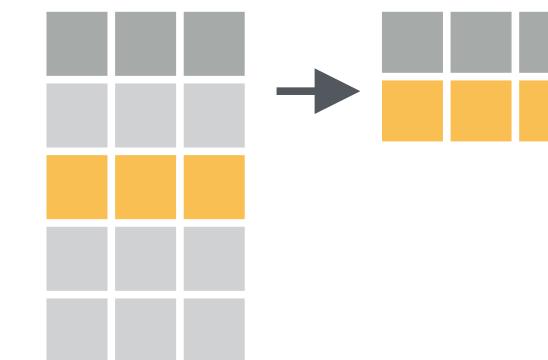


The penguins dataset

Snippet of the dataset:

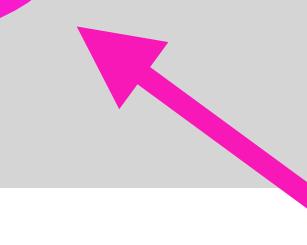
species	island	bill_len	bill_dep	flipper_len	body_mass	sex	year
Adelie	Torgersen	39.1	18.7	181	3750	male	2007
Adelie	Torgersen	39.5	17.4	186	3800	female	2007
Adelie	Torgersen	40.3	18.0	195	3250	female	2007
Adelie	Torgersen	NA	NA	NA	NA	NA	2007
Adelie	Torgersen	36.7	19.3	193	3450	female	2007

filter()



filter observations by their values

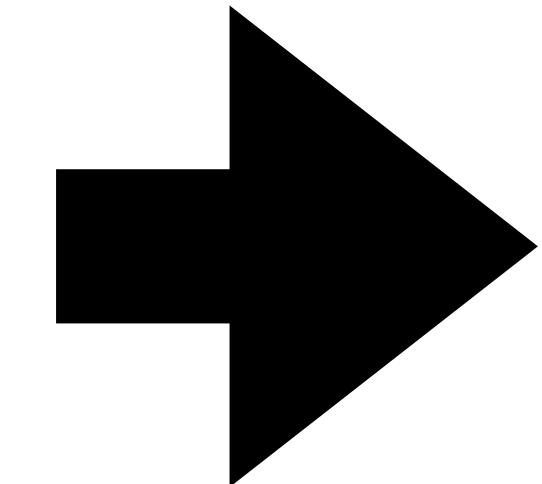
```
gentoo <- penguins %>%  
  filter(species == "Gentoo")
```



two equals signs for “equals to”

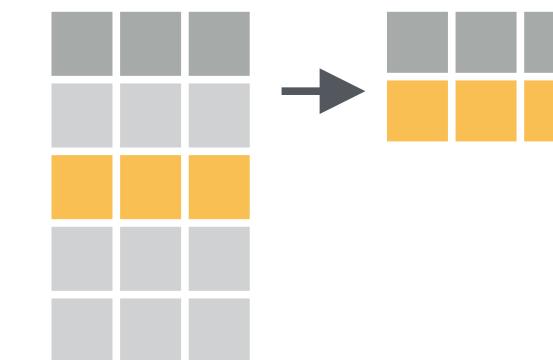
Goal:

full
penguins
dataset



only
gentoo
penguins

filter()

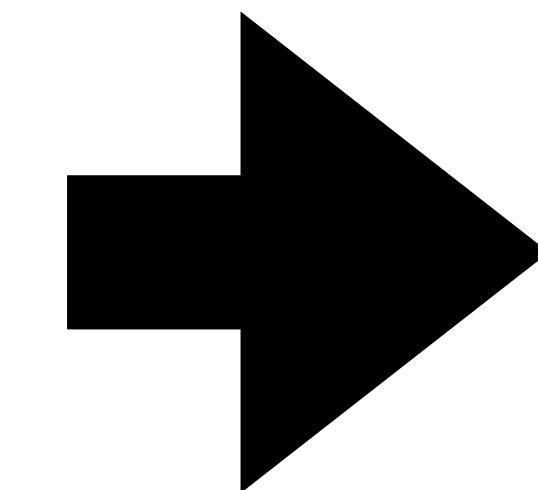


filter observations by their values

```
big_penguins <- penguins %>%  
  filter(body_mass > 4500)
```

Goal:

full
penguins
dataset

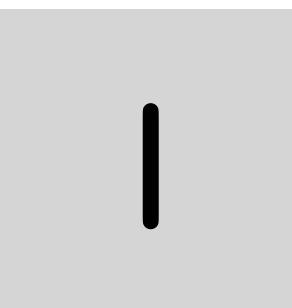


body mass
of more
than 4500

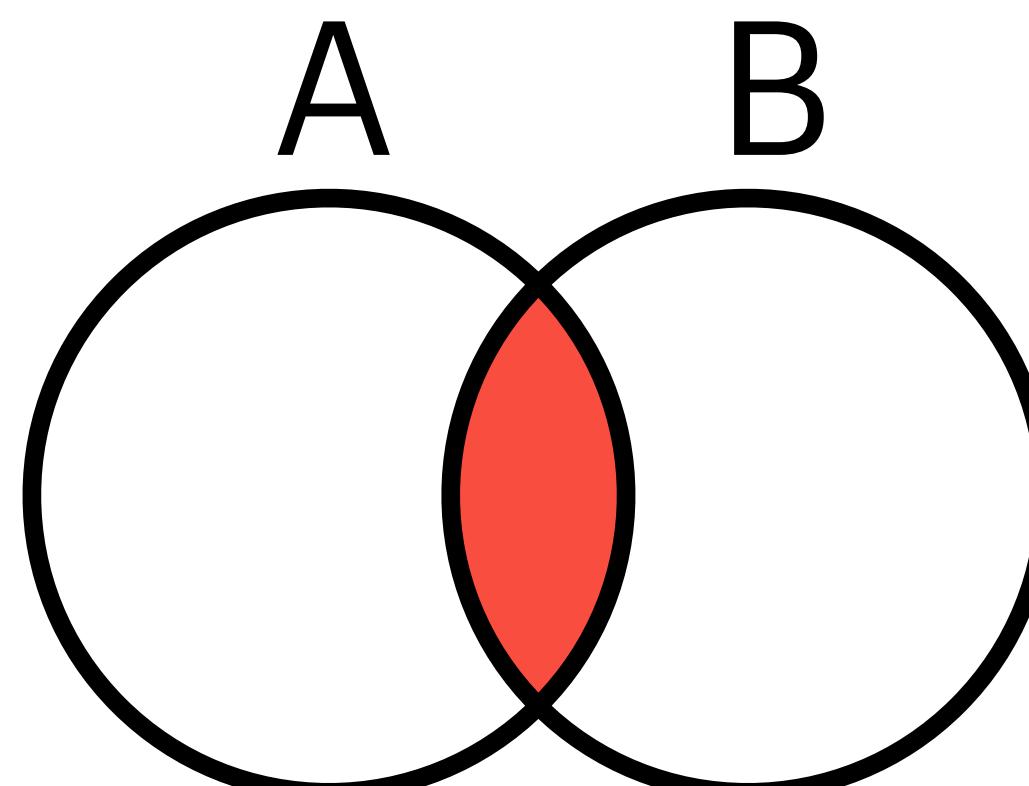
Comparison “operators” in R

Operator	Description
>	Greater than
<	Less than
==	Equals to
!=	Not equal to
>=	Greater than or equal to
<=	Less than or equal to

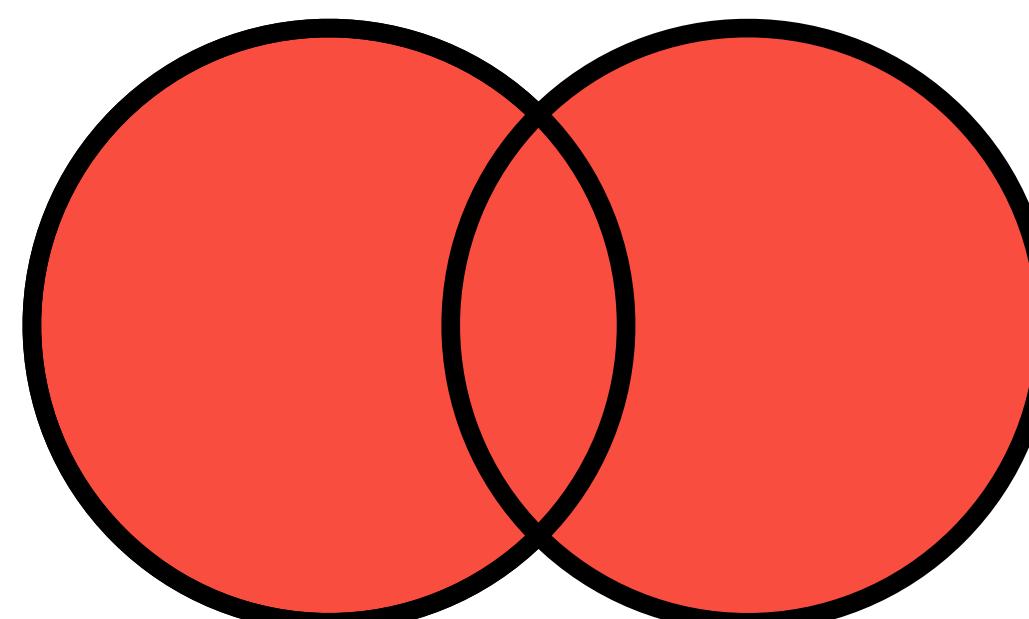
Combine multiple statements

 & “and”
 | “or”

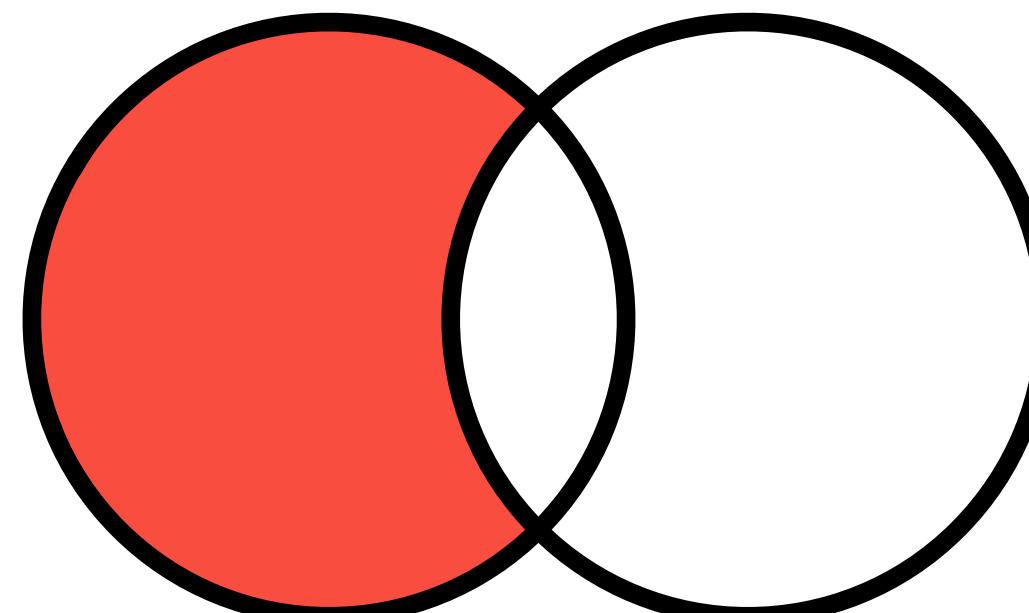
Comparison “operators” in R



A and B

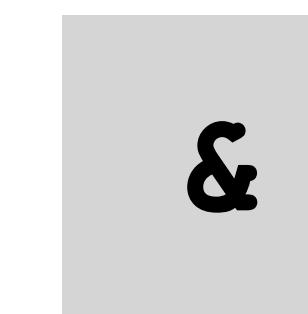


A or B

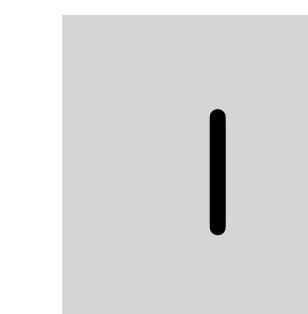


A and not B

Combine multiple statements

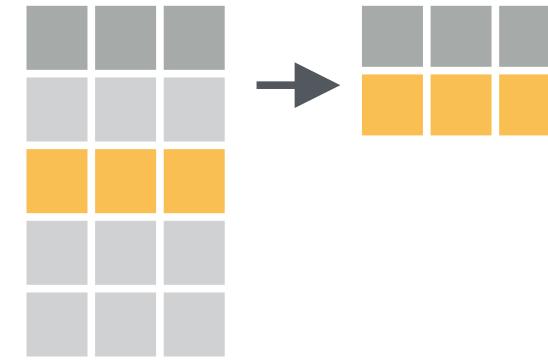


“and”



“or”

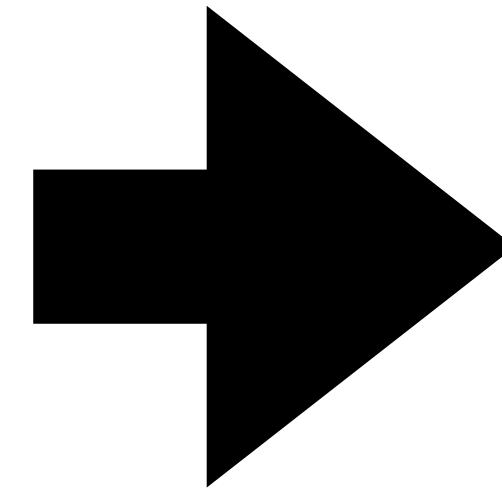
filter()



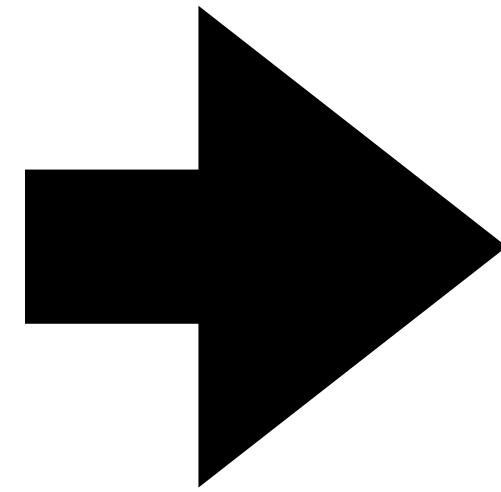
filter observations by their values

Goal:

full
penguins
dataset

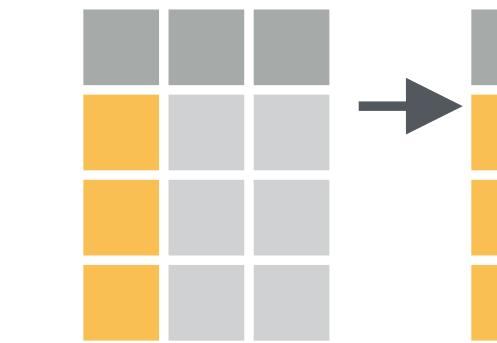


body mass
of more
than 4500



only
gentoo
penguins

select()

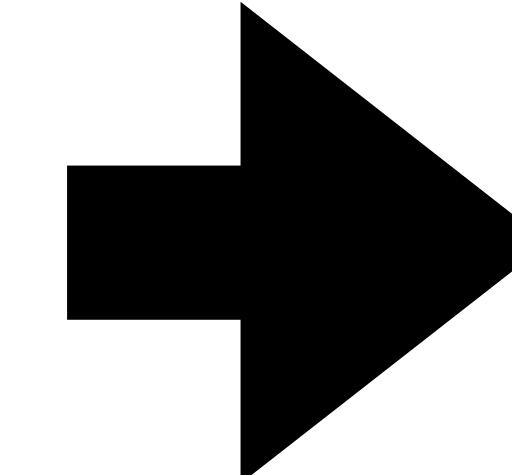


select variables by their column names

```
new_penguins <- penguins %>%  
  select(species, sex, body_mass)
```

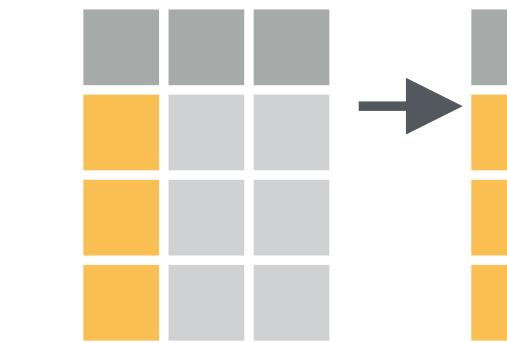
Goal:

full
penguins
dataset



dataset with
species,
sex, and
body_mass

select()

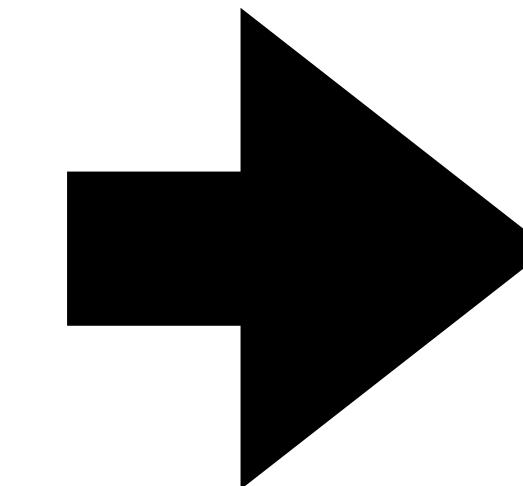


select variables by their column names

```
numeric_penguins <- penguins %>%  
  select(where(is.numeric))
```

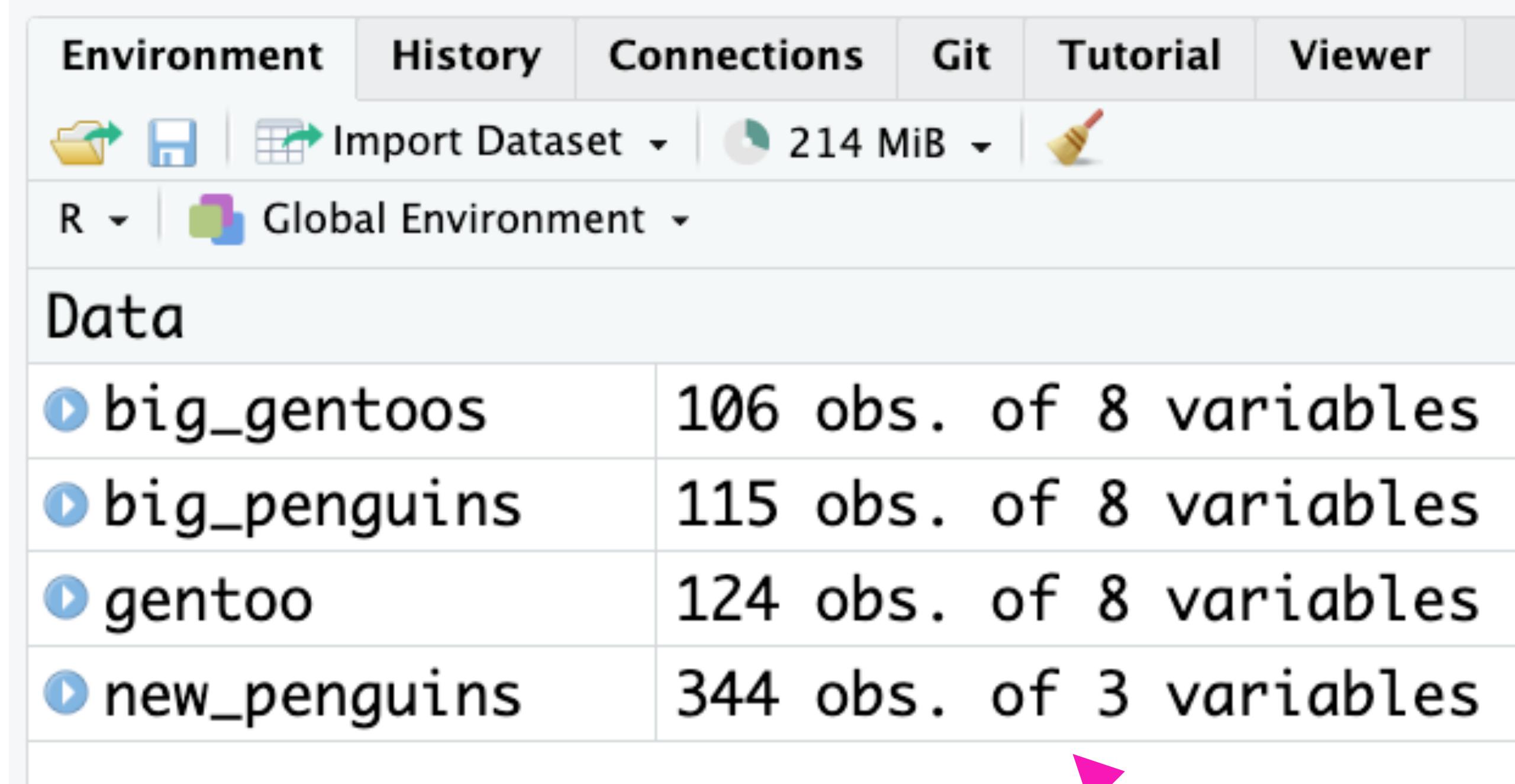
Goal:

full
penguins
dataset



dataset with
numeric
columns

How does your R environment look now?



The screenshot shows the RStudio interface with the 'Environment' tab selected. The global environment contains four dataframes:

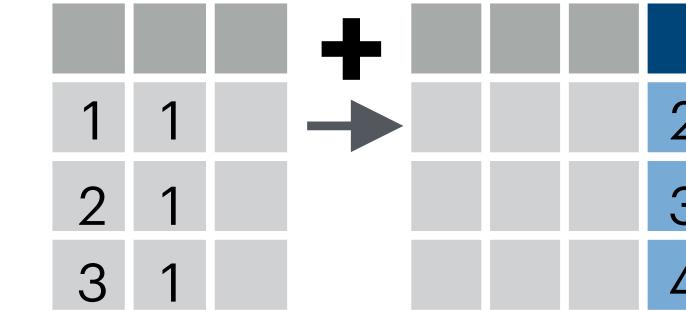
Object	Description
big_gentoos	106 obs. of 8 variables
big_penguins	115 obs. of 8 variables
gentoo	124 obs. of 8 variables
new_penguins	344 obs. of 3 variables

A pink curly brace on the left groups the first three dataframes, labeled "more dataframes". A pink arrow points from the text "all of different dimensions" to the bottom-right dataframe, "new_penguins".

more dataframes {

all of different dimensions

mutate()



mutate (or change) existing variables to create new ones

add to our original dataset

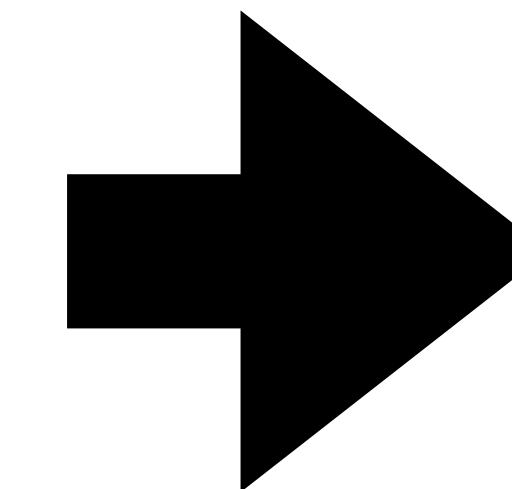
```
penguins <- penguins %>%  
  mutate(bill_shape) = bill_len/bill_dep)
```

new variable name

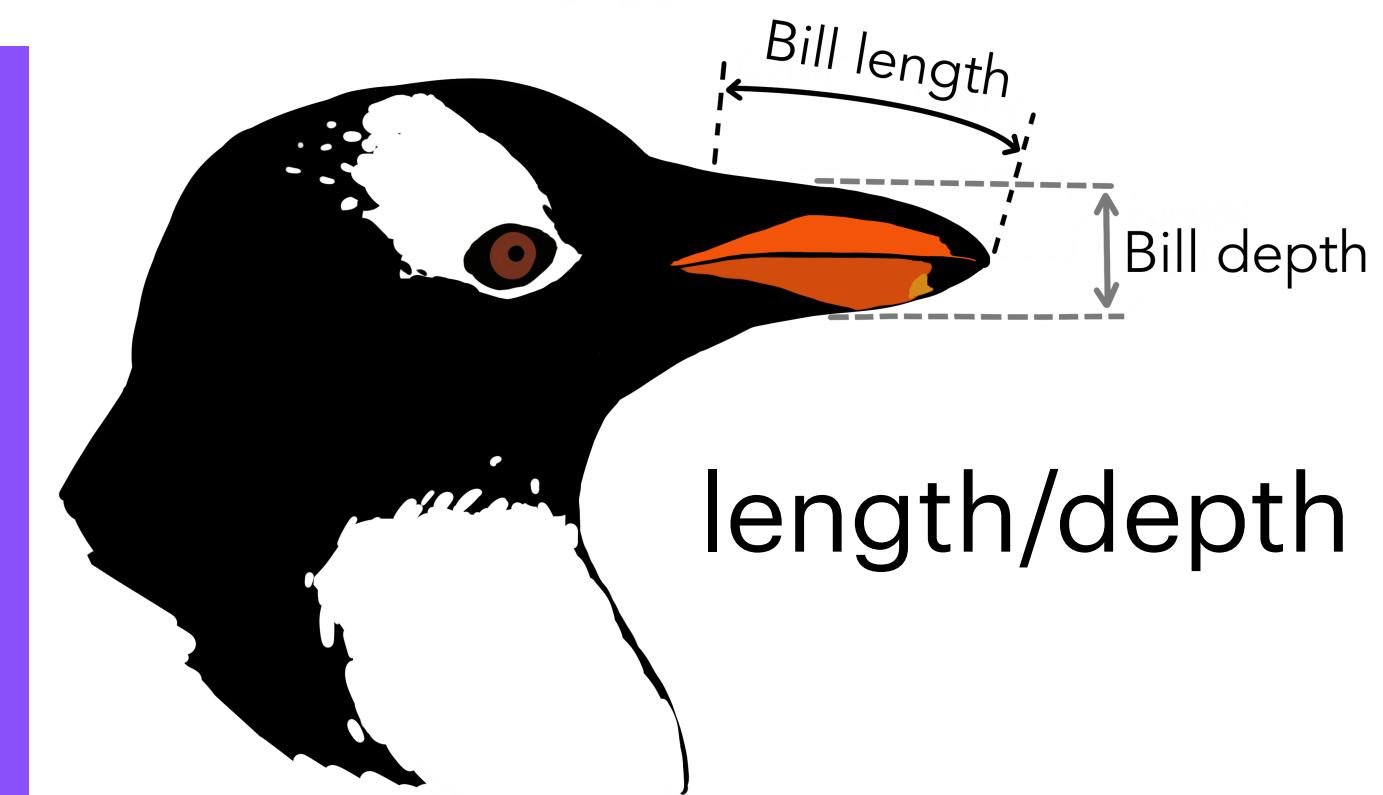
the calculation

Goal:

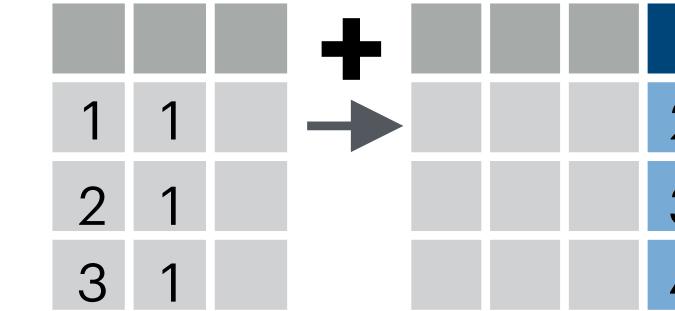
full
penguins
dataset



new
variable for
bill shape



mutate()



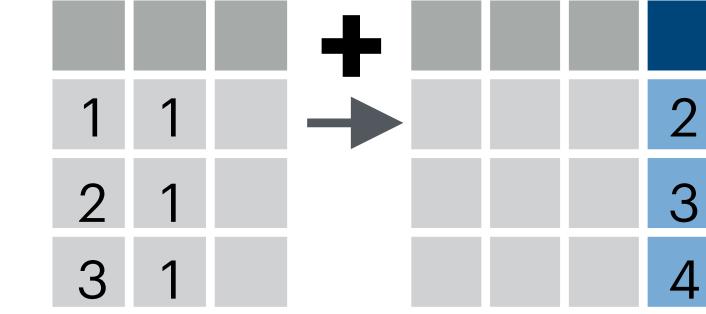
mutate (or change) existing variables to create new ones

```
df <- df %>%  
  mutate(<NEW_VARIABLE> = <EQUATION>)
```

```
df_new <- df %>%  
  mutate(<NEW_VARIABLE> = <EQUATION>)
```

instead create a new data frame

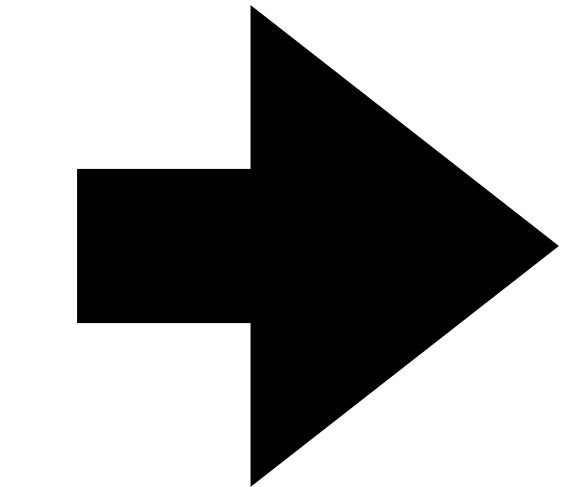
mutate()



mutate (or change) existing variables to create new ones

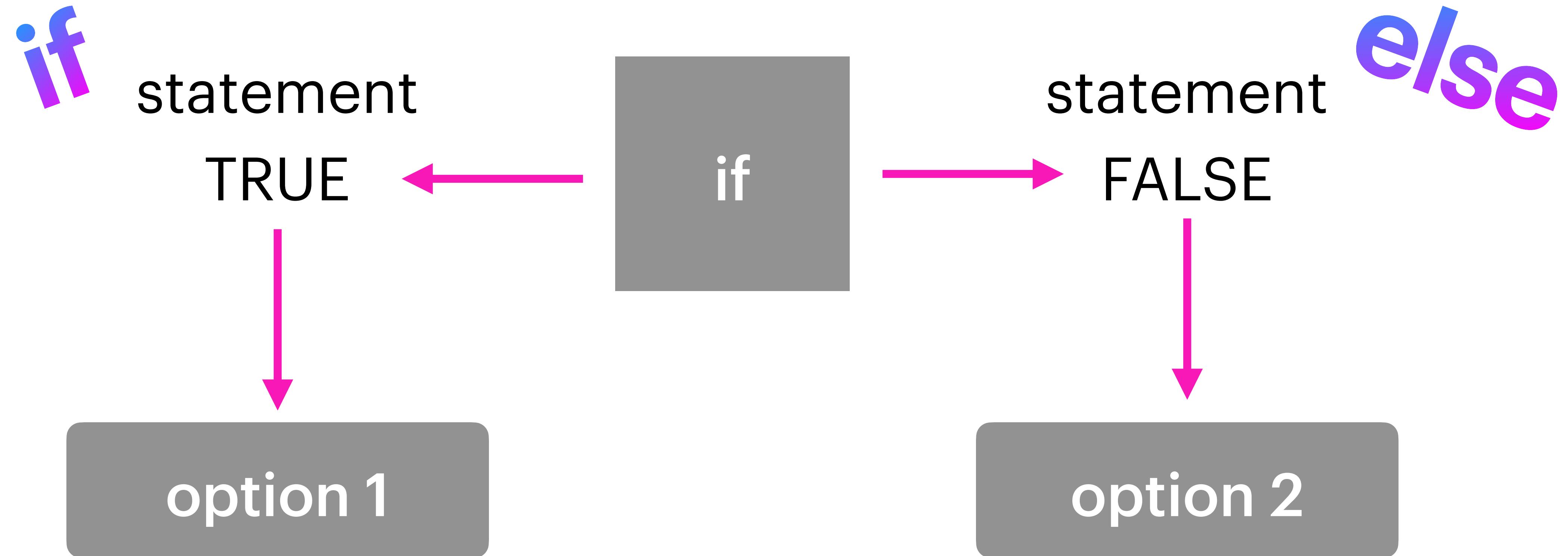
Goal:

full
penguins
dataset



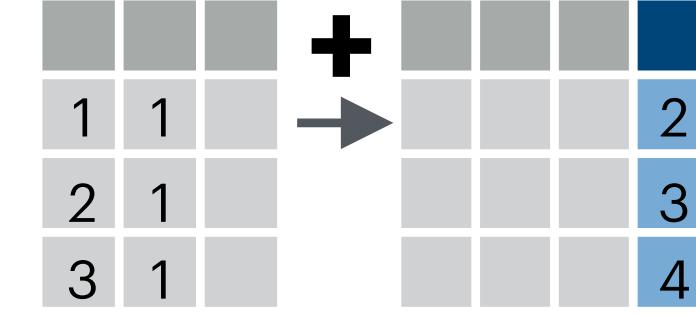
new
variable
year times
body mass

ifelse() function



```
ifelse(<test>, <true option>, <false option>)
```

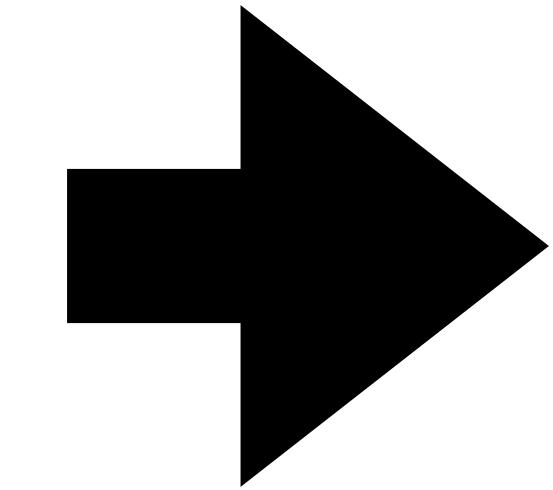
mutate()



mutate (or change) existing variables to create new ones

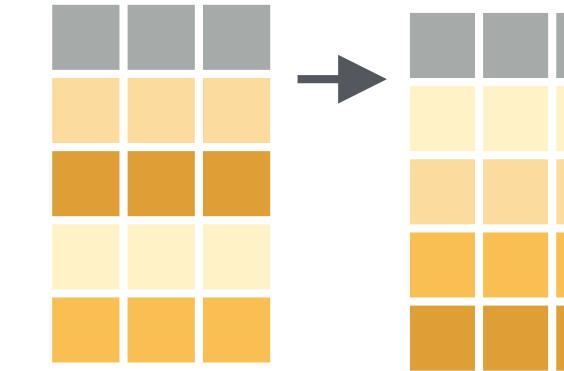
Goal:

full
penguins
dataset



new
variable
“Adelie”
(yes or no)

arrange ()

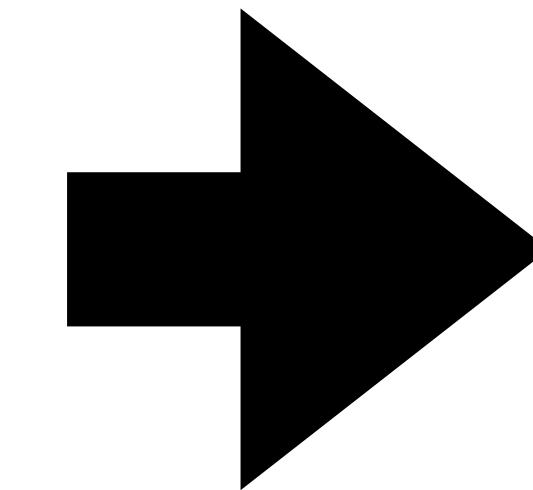


arrange observations by their values

```
penguins <- penguins %>%  
  arrange(island)
```

Goal:

full
penguins
dataset



dataset
arranged
by island

arrange ()



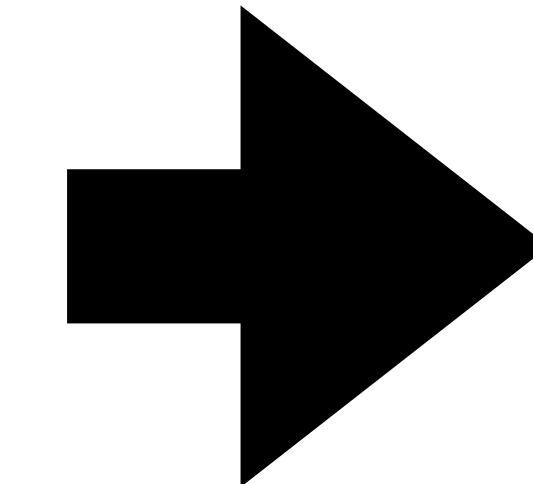
arrange observations by their values

```
penguins <- penguins %>%  
  arrange(desc(body_mass))
```

“descending”

Goal:

full
penguins
dataset



dataset
arranged by
mass (largest
to smallest)

Make nice tables in Quarto

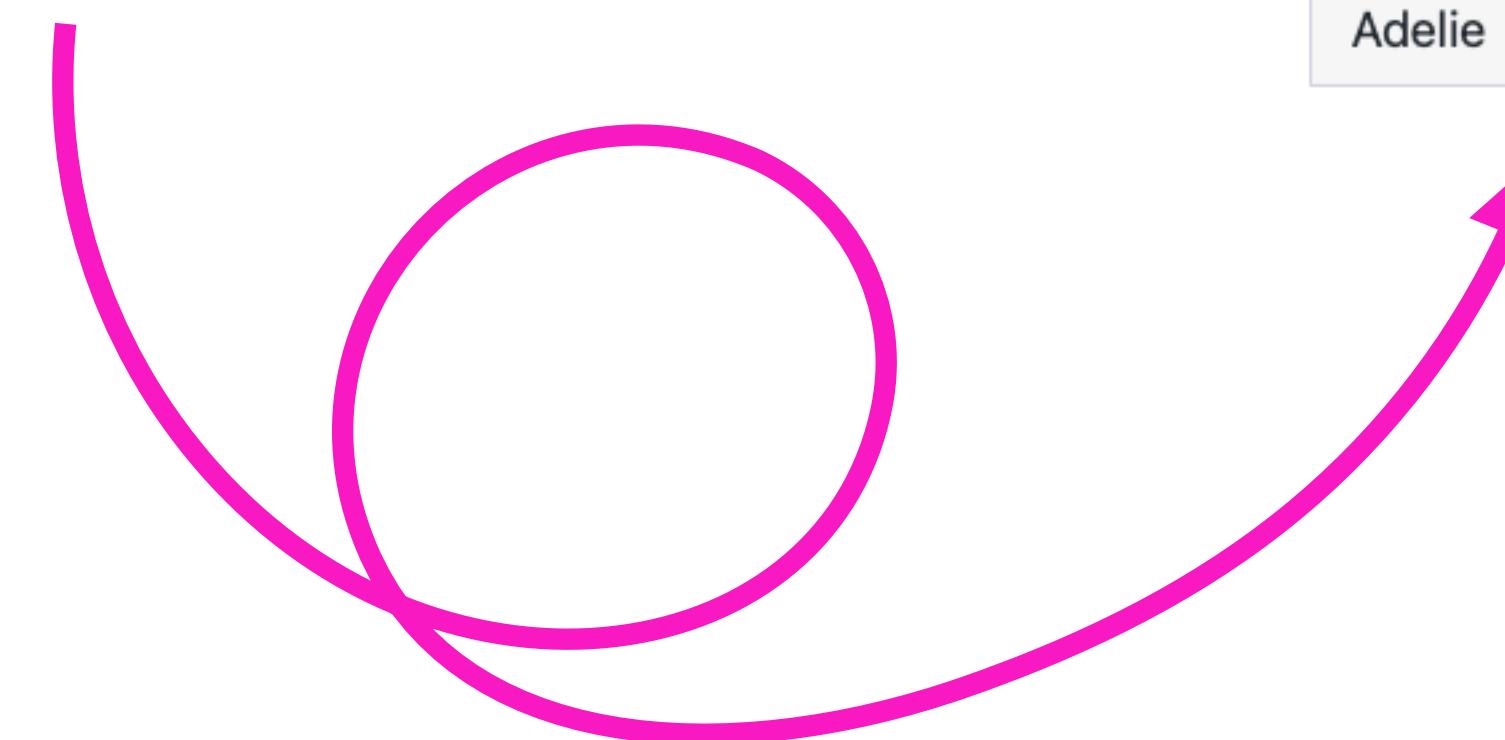
Load (or install + load) the R Package `knitr` to use the `kable()` function

`head(penguins)`

	species	island	bill_len	bill_dep	flipper_len	body_mass	sex	year
1	Adelie	Torgersen	39.1	18.7	181	3750	male	2007
2	Adelie	Torgersen	39.5	17.4	186	3800	female	2007
3	Adelie	Torgersen	40.3	18.0	195	3250	female	2007
4	Adelie	Torgersen	NA	NA	NA	NA	<NA>	2007
5	Adelie	Torgersen	36.7	19.3	193	3450	female	2007
6	Adelie	Torgersen	39.3	20.6	190	3650	male	2007

`kable(head(penguins))`

species	island	bill_len	bill_dep	flipper_len	body_mass	sex	year
Adelie	Torgersen	39.1	18.7	181	3750	male	2007
Adelie	Torgersen	39.5	17.4	186	3800	female	2007
Adelie	Torgersen	40.3	18.0	195	3250	female	2007
Adelie	Torgersen	NA	NA	NA	NA	NA	2007
Adelie	Torgersen	36.7	19.3	193	3450	female	2007
Adelie	Torgersen	39.3	20.6	190	3650	male	2007



Further reading

- Data Transformation with `dplyr` [Cheatsheet](#)
- Boehmke 2016, [Data Wrangling with R](#)
- Alston & Rick 2021. [A Beginner's Guide to Conducting Reproducible Research. The Bulletin of the Ecological Society of America](#)